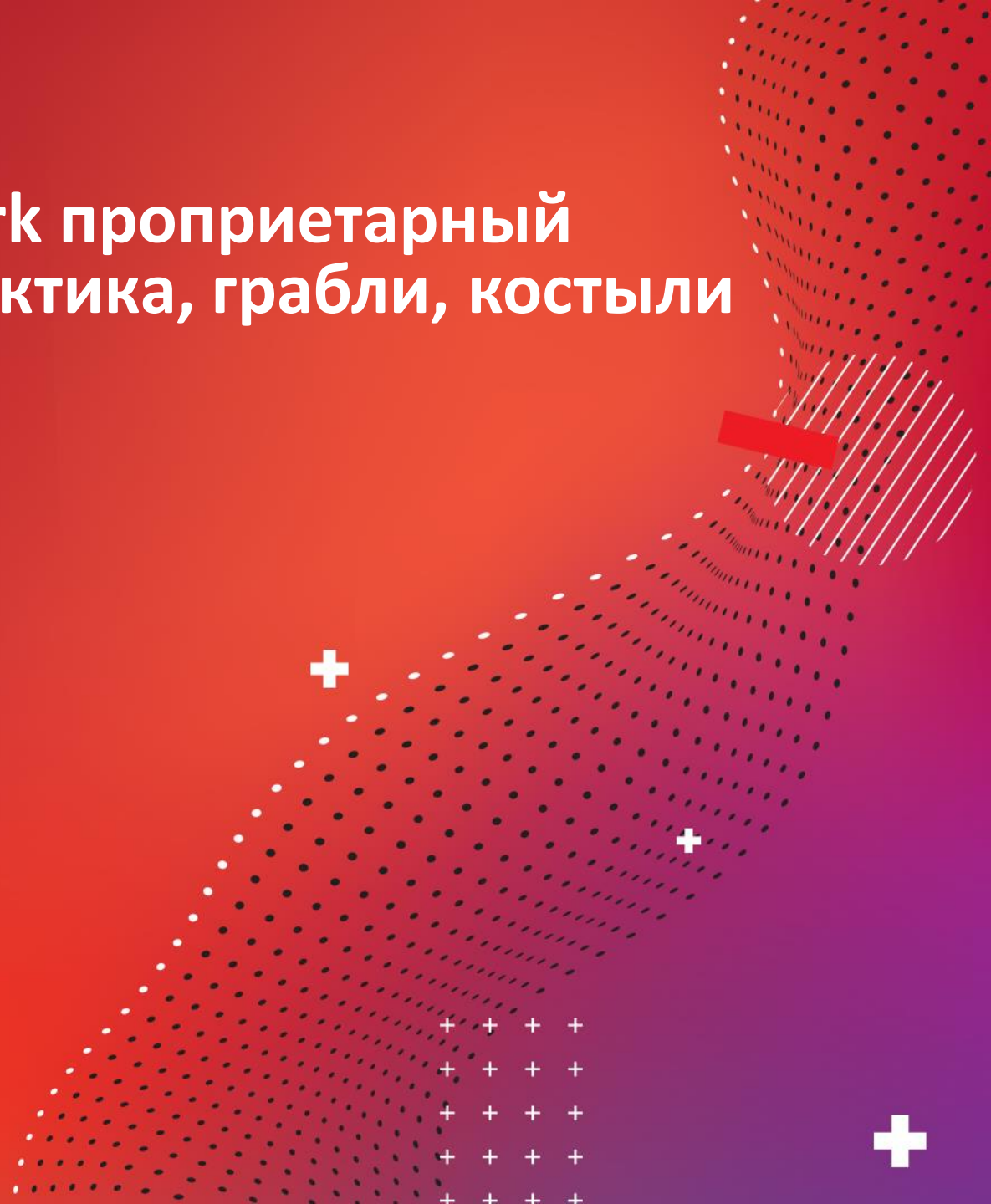


# Как подключить к Apache Spark проприетарный источник данных: теория, практика, грабли, костыли

Александра Белоусова  
Яндекс.Go



**HighLoad++**  
Весна 2021



# О чём я буду говорить

- Мотивация и постановка задачи
- Быстрое и простое решение
- Сложное производительное решение
- Грабли и костыли

Что надо было сделать?

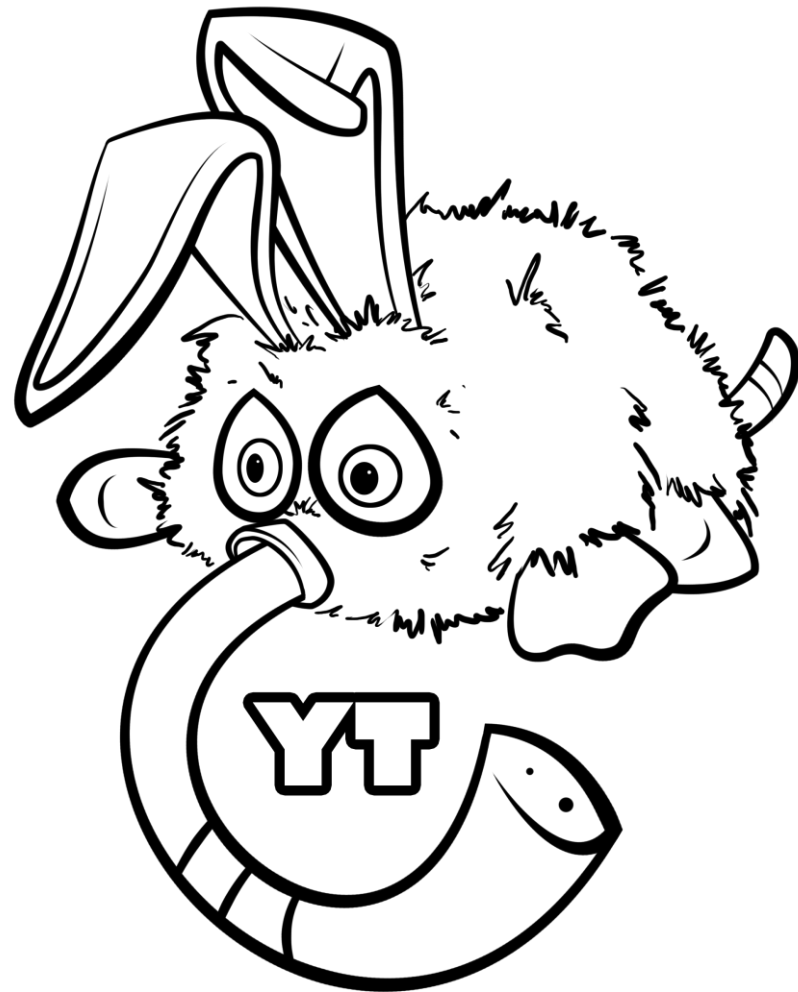
# YT

- Data Storage
- Scheduler
- Map Reduce



# YT

- Миллион ядер
- Эксбибайт (EiB) диска
- Около 10 тысяч пользователей
- И это только на одном кластере!



# Spark + YT

- Запустить на ресурсах кластера YT
  - Читать данные из YT
  - Писать данные в YT
- 
- Показать wall time и cpi time лучше, чем у существующих решений

# YT – файловая система

## HDFS

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> dataset_1
|           |-----> part-0.parquet
|           |-----> part-1.parquet
|           |-----> part-2.parquet
|       |---> dataset_2
|           |-----> part-0.parquet
|           |-----> part-1.parquet
```

## YT

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |-----> example_table_1
|       |-----> example_table_2
```

# Таблица в YT

## User level

#	a	b	c
0	"cat"	true	1
1	"dog"	true	20
2	"otter"	false	350

### Attributes:

schema  
a - String  
b - Boolean  
c - Integer

rowCount

## Storage level

### Columnar chunk 1

rows 0 - 200

### Columnar chunk 2

rows 201 - 400

### Columnar chunk 3

rows 401 - 600

### Attributes:

schema  
a - String  
b - Boolean  
c - Integer

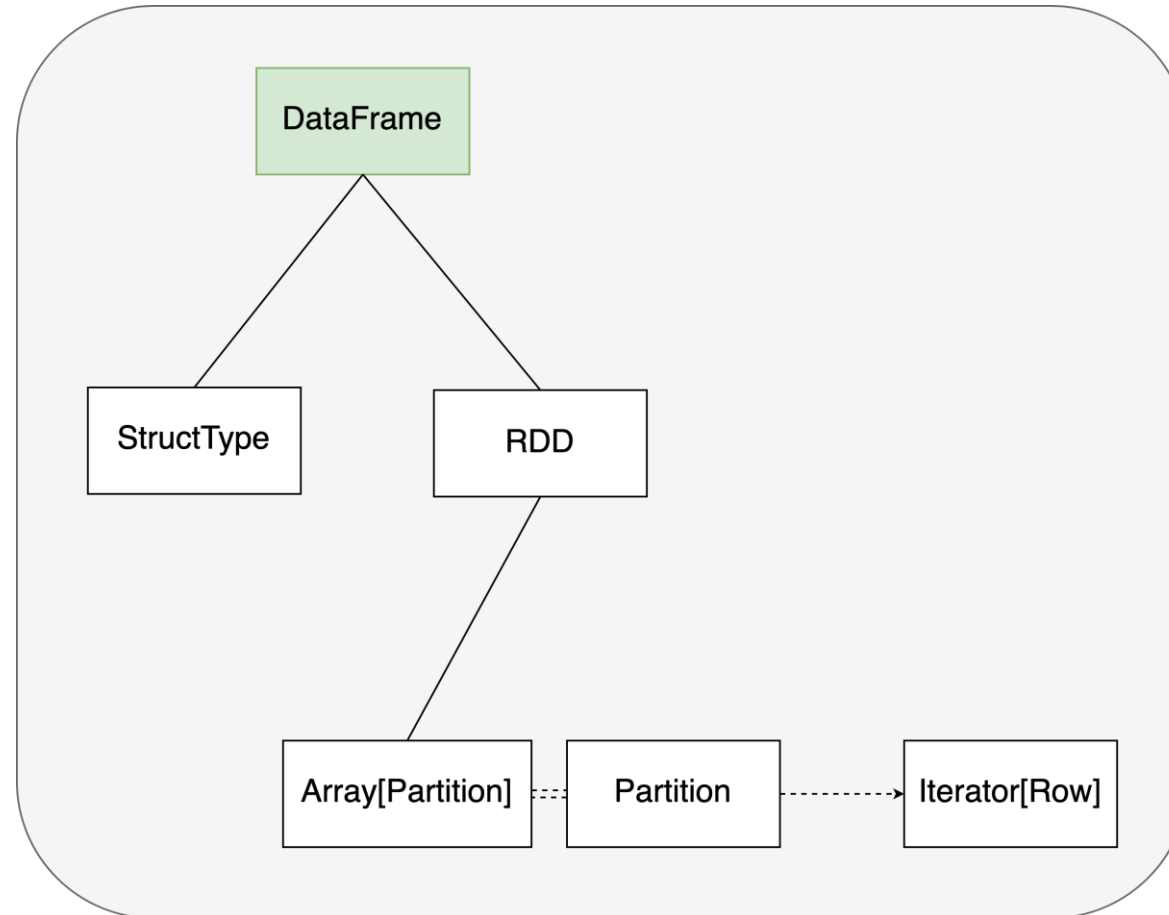
rowCount

Простое решение: чтение

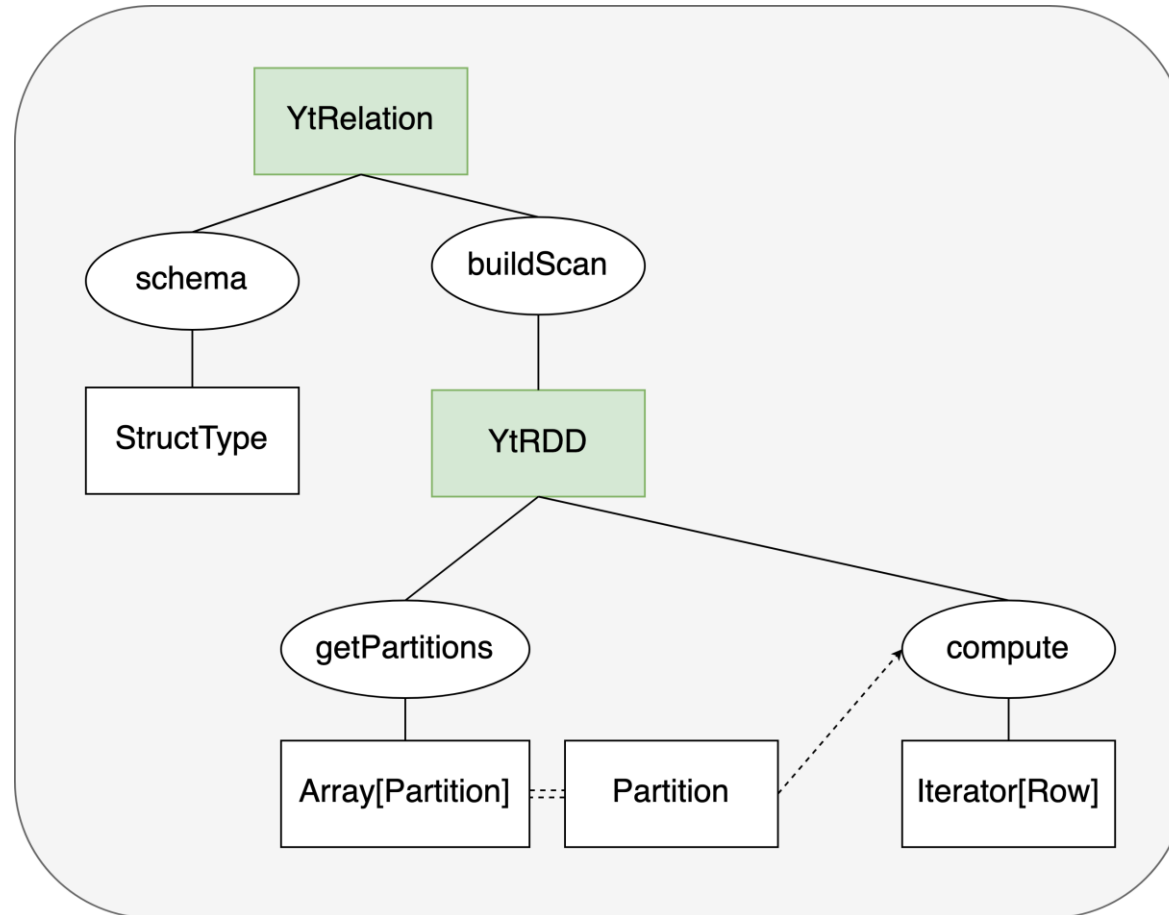
# Дисклеймер

- API – набор классов / интерфейсов / методов, которые мы *вызываем*
- SPI – набор классов / интерфейсов / методов, которые мы *расширяем или реализуем*

# Spark SPI

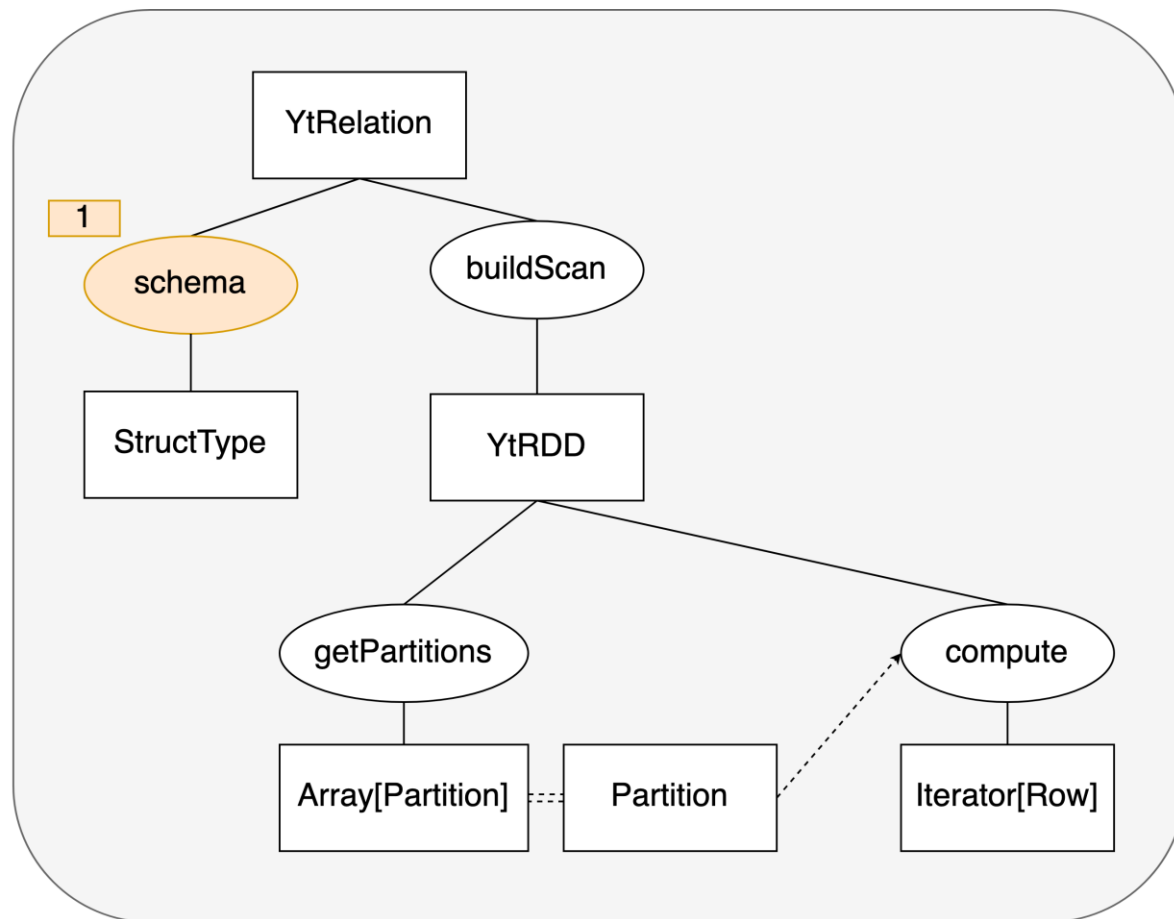


# Spark SPI

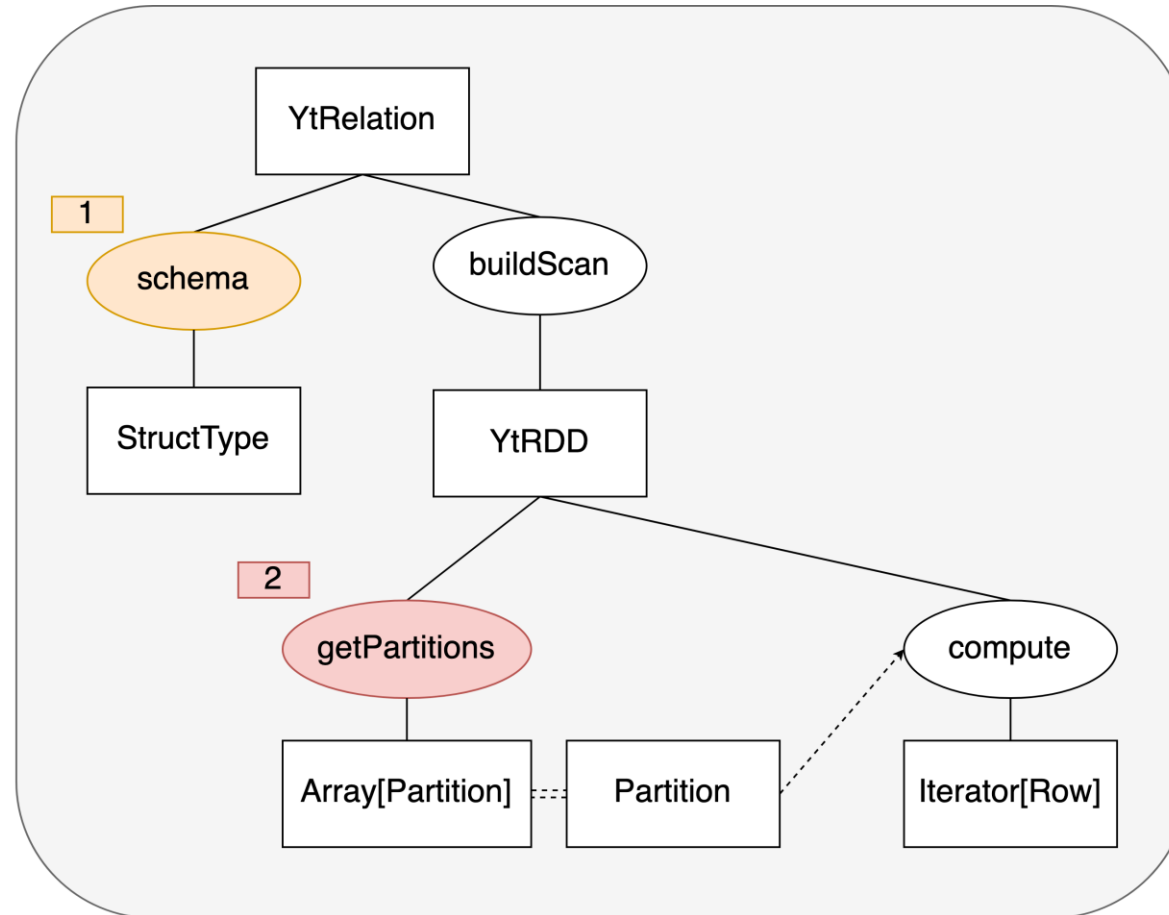




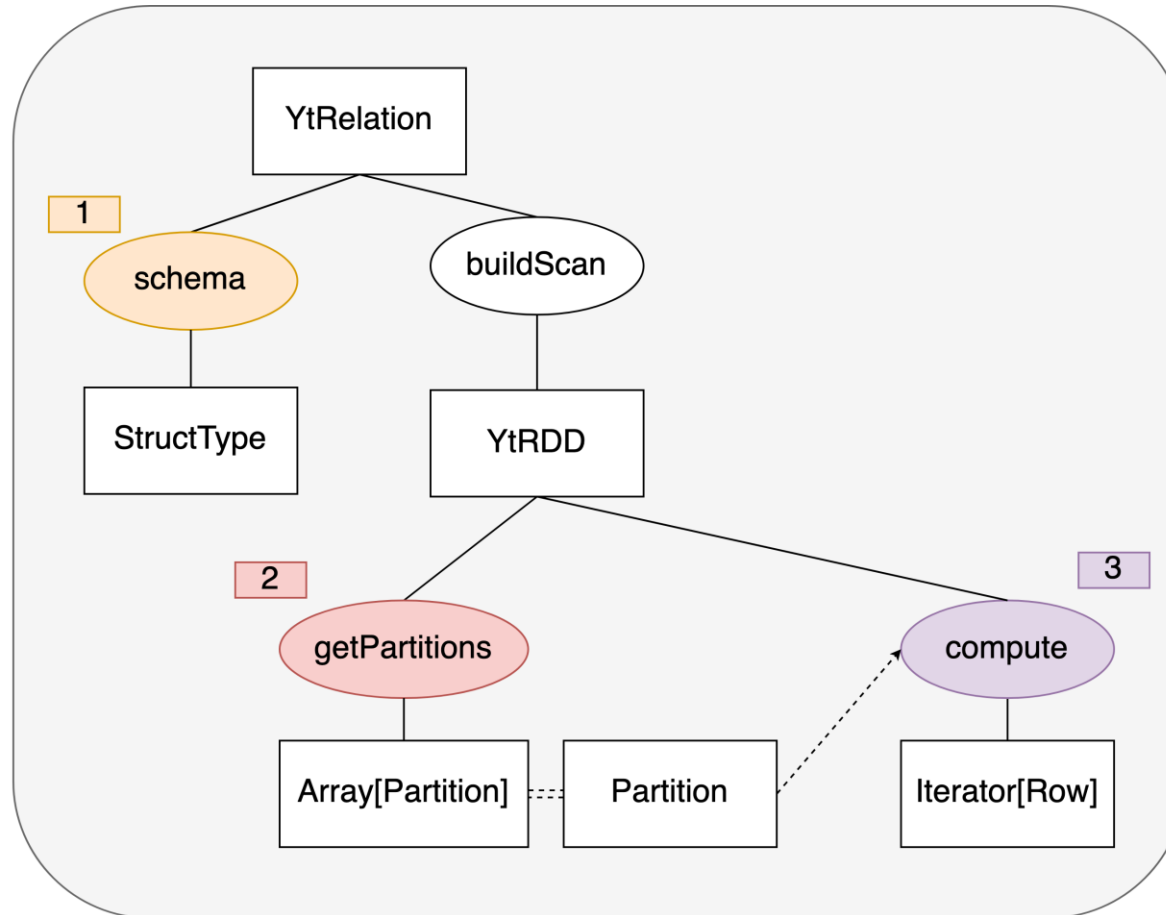
# Получение схемы



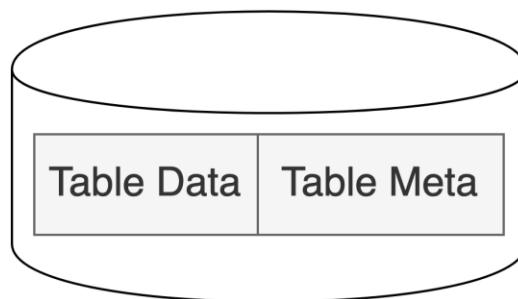
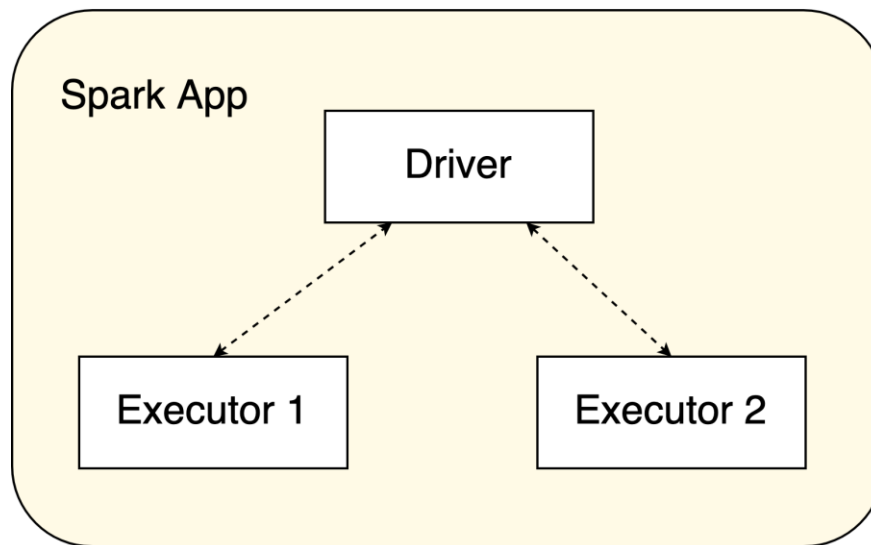
# Получение партиций



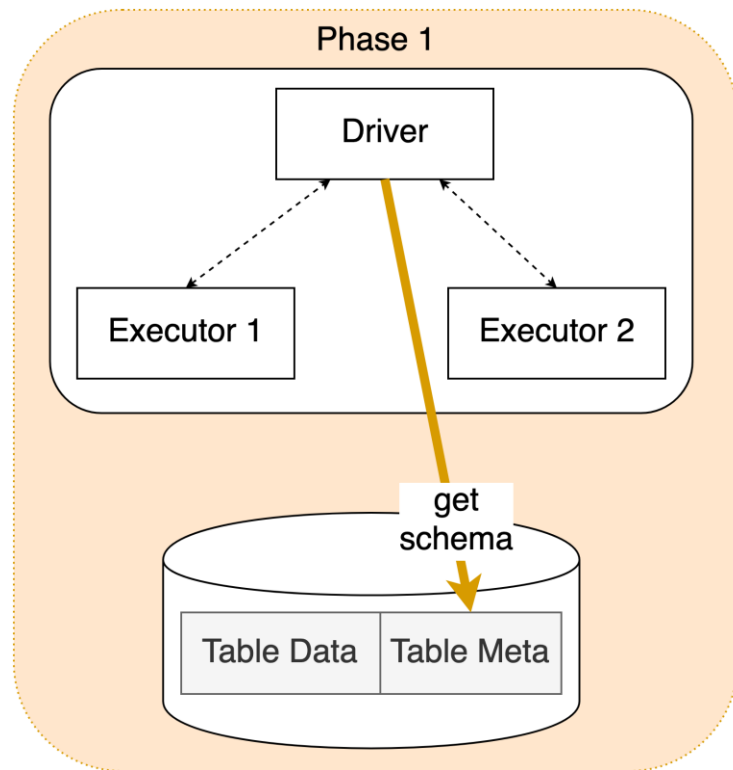
# Чтение партиции



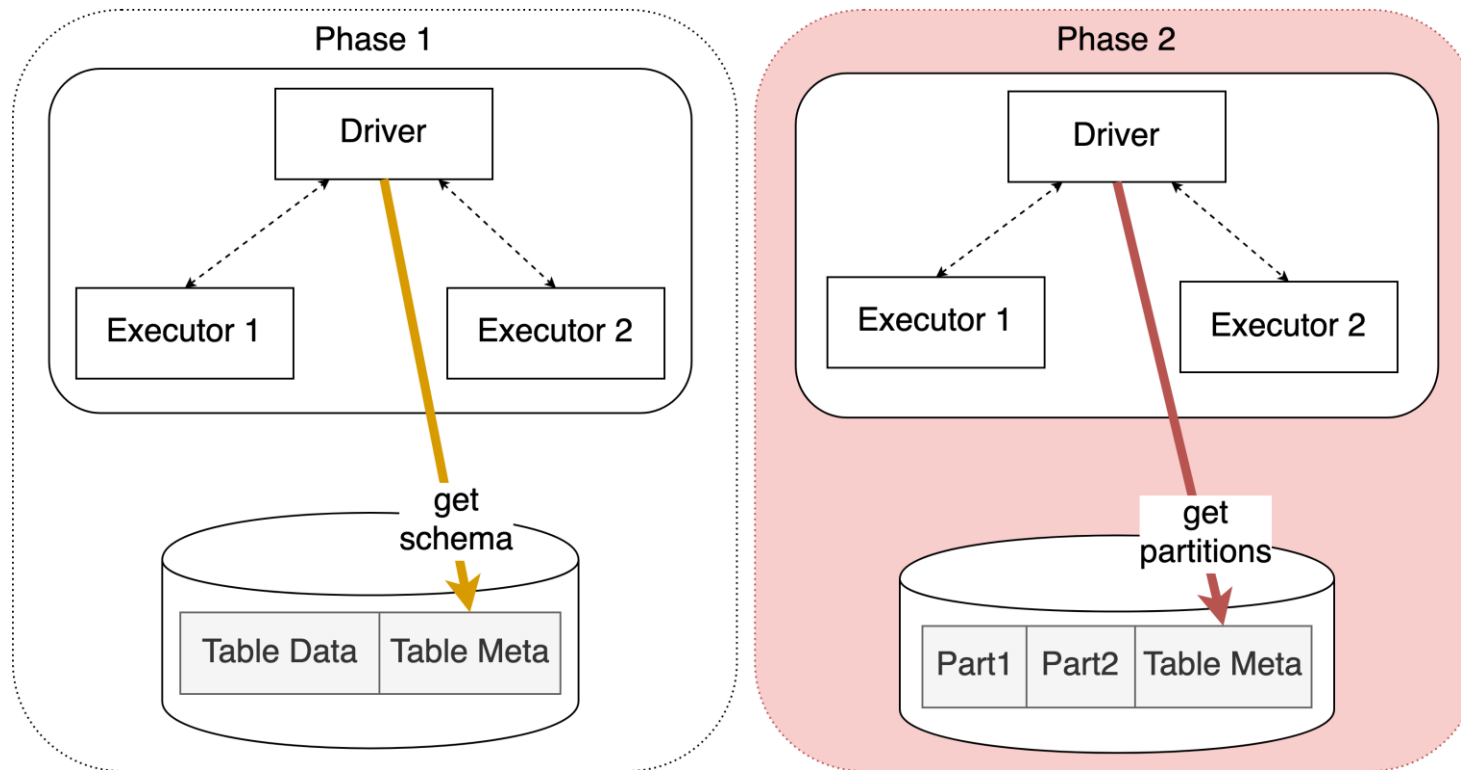
# Как это выполняется



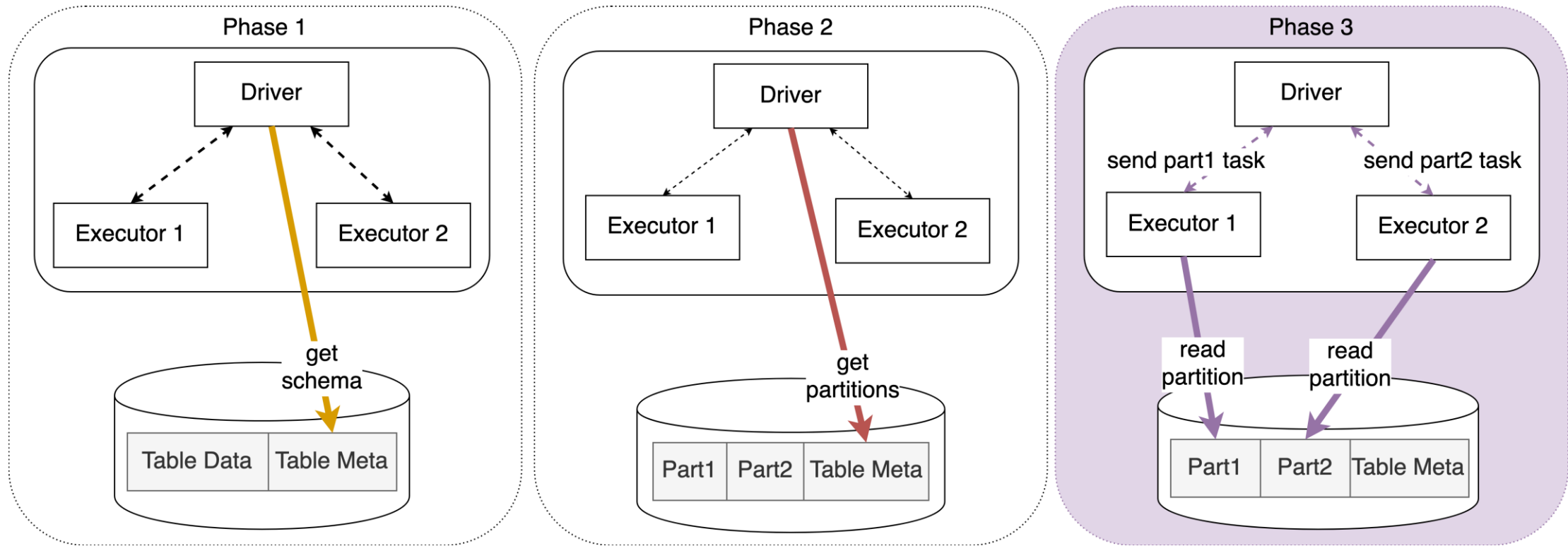
# Получение схемы



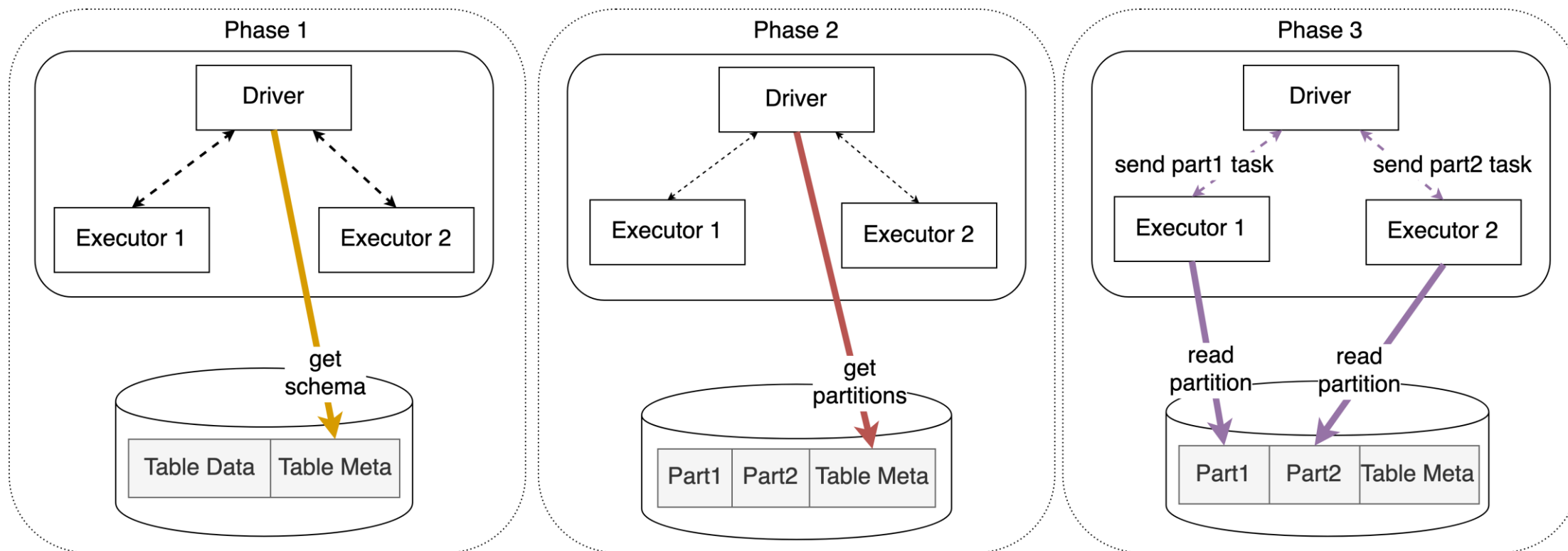
# Получение партиций



# Чтение партиции

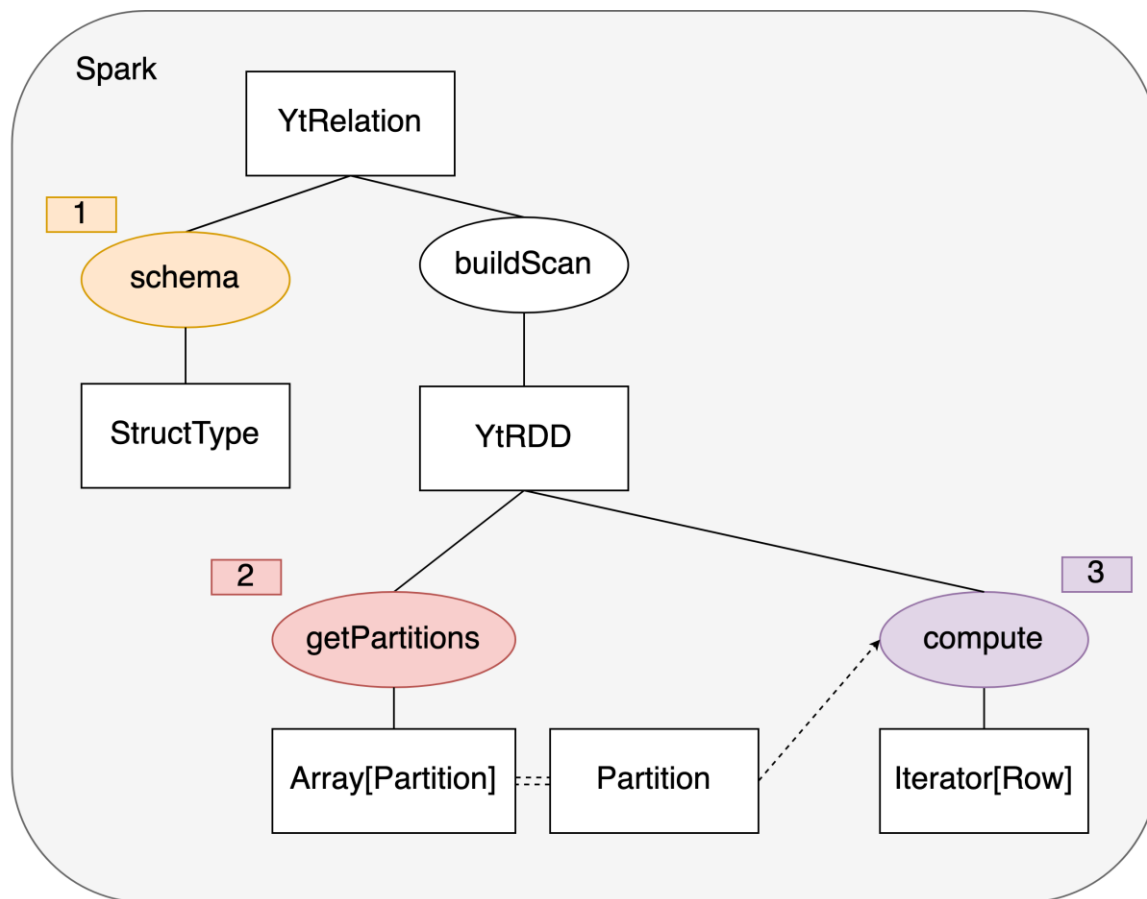


# Как это выполняется

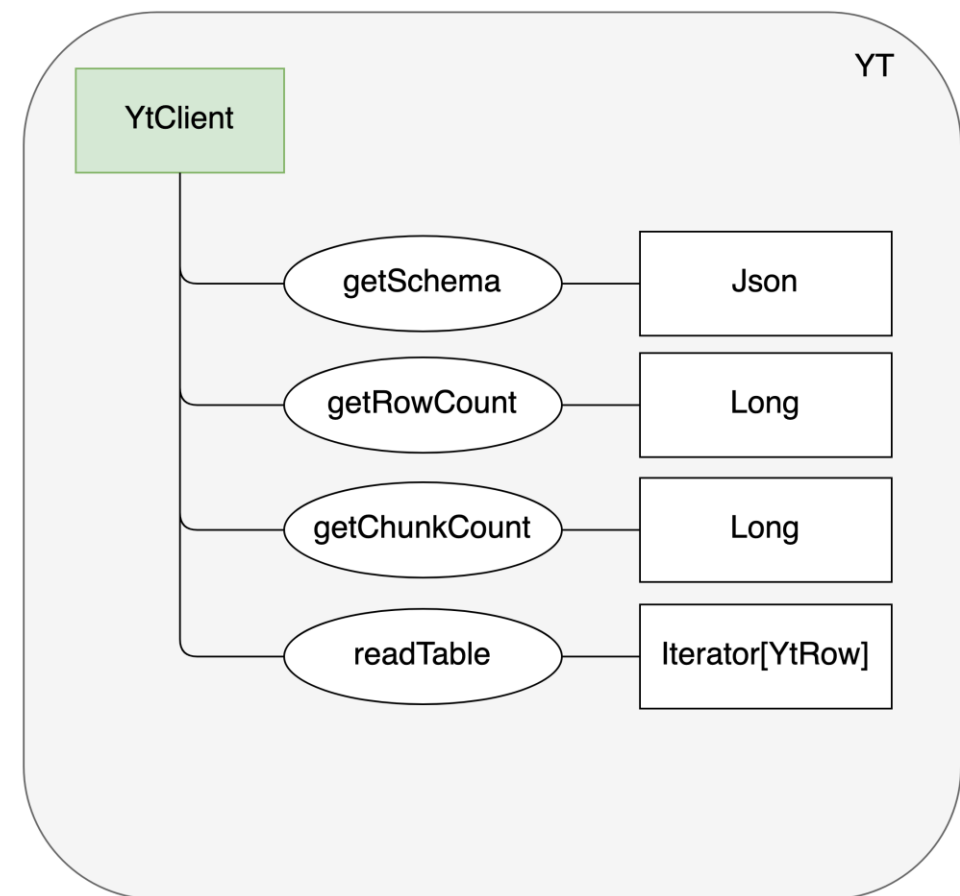
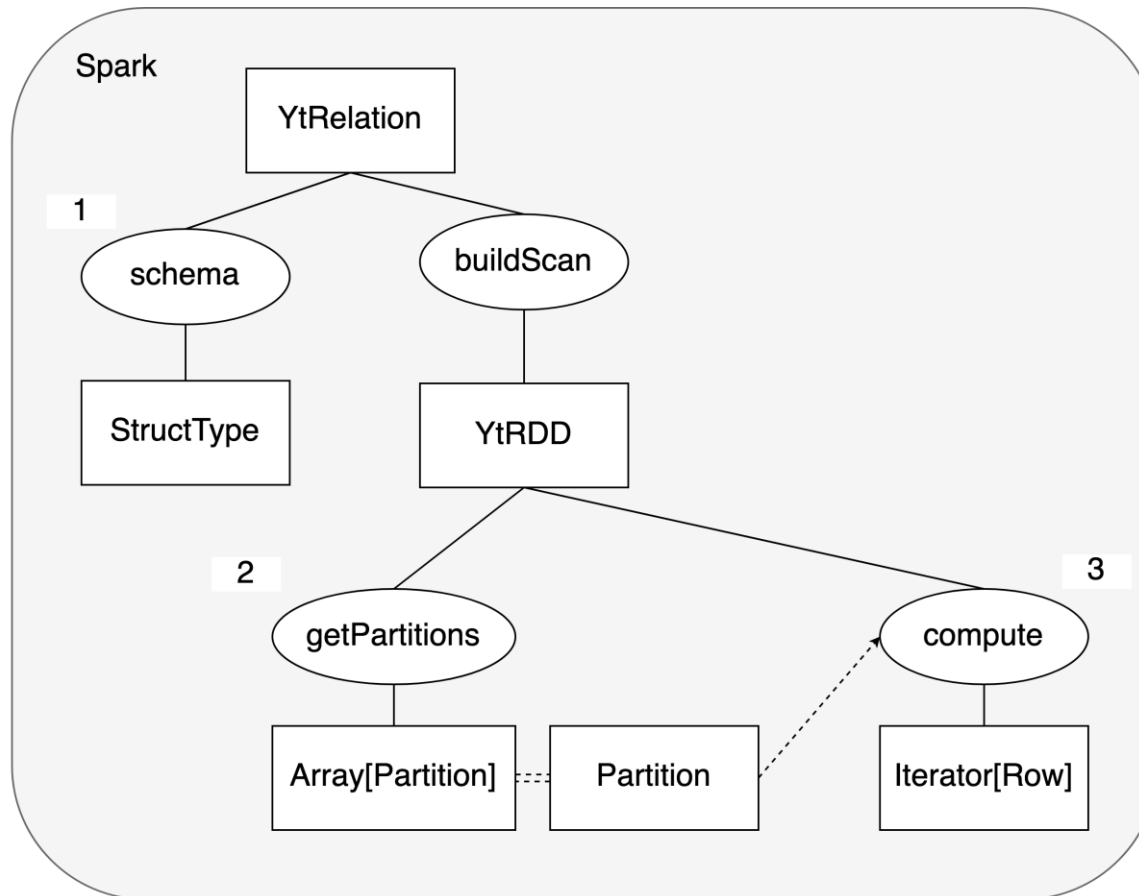




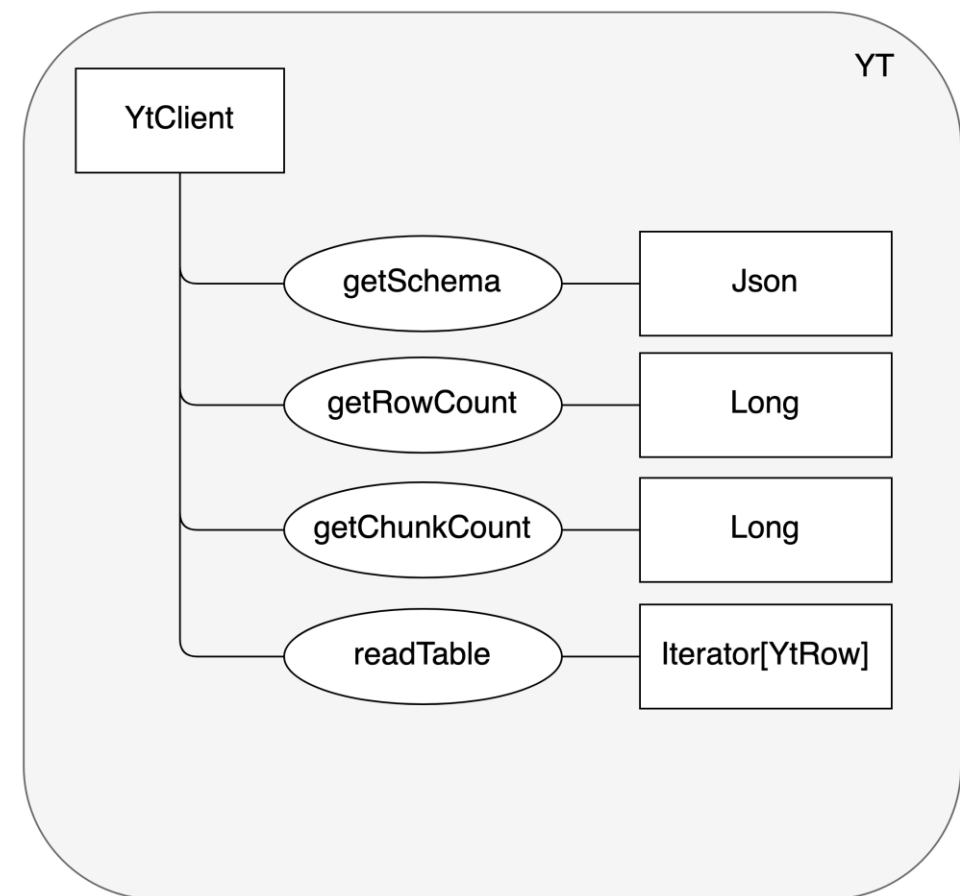
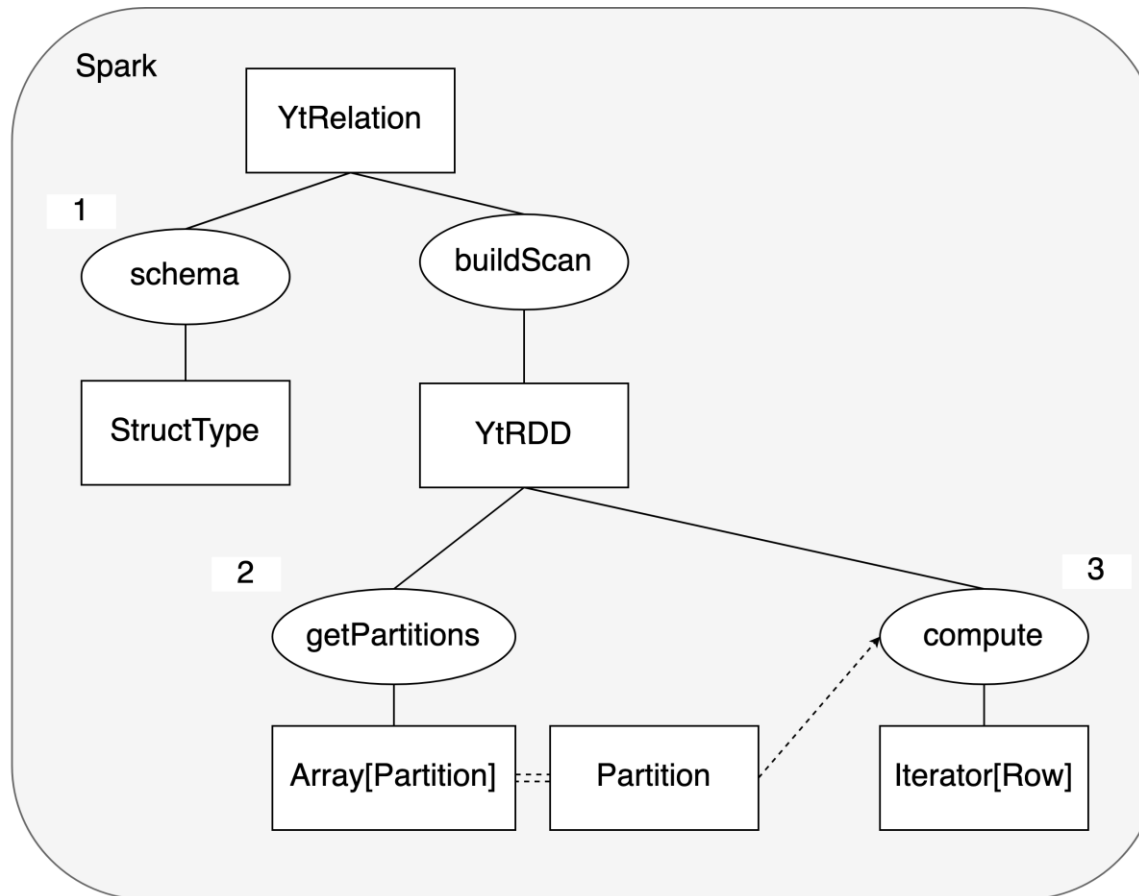
# Как это выглядит в коде



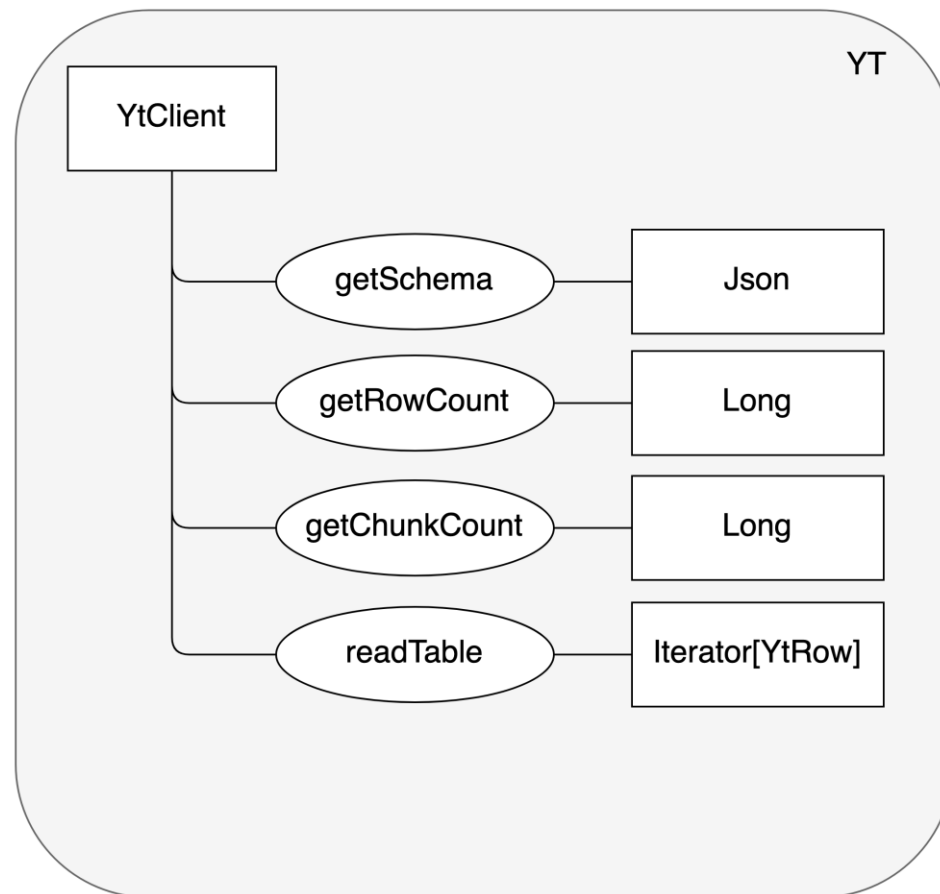
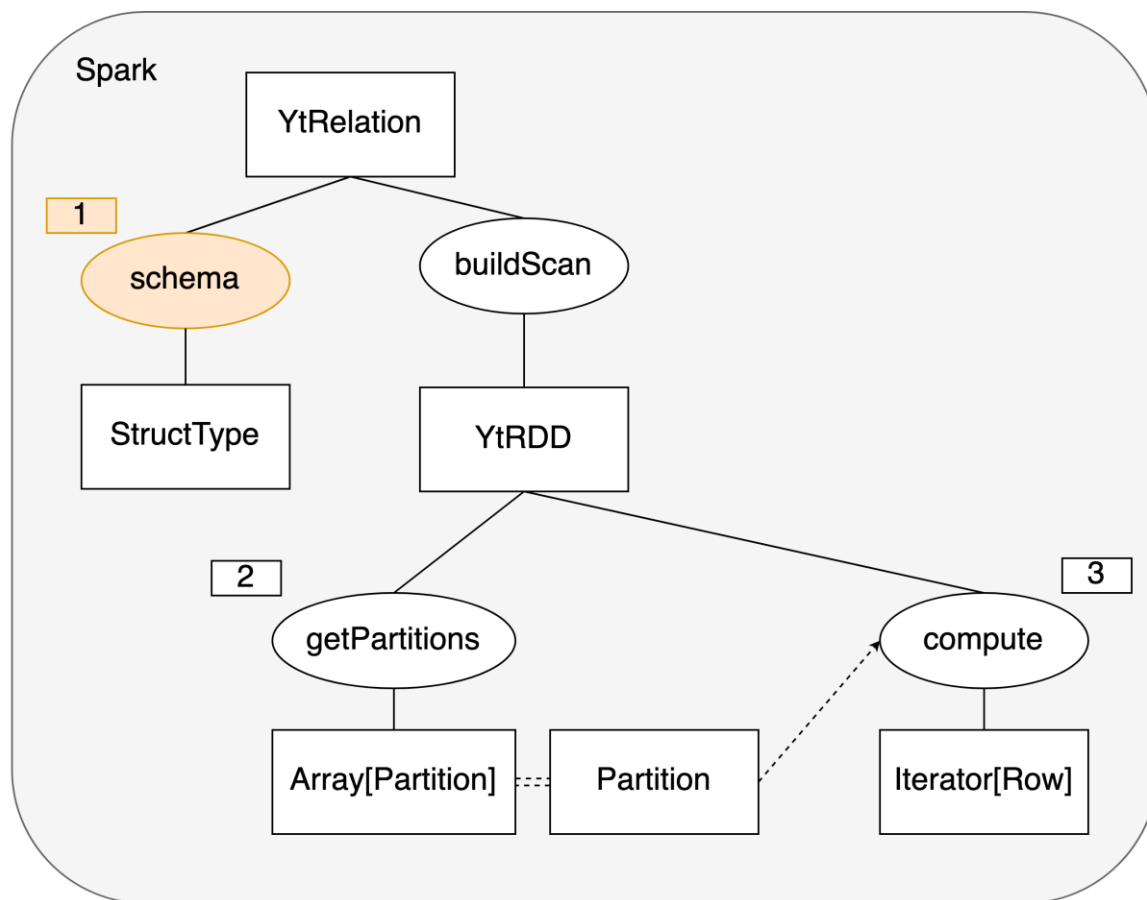
# YT API



# SPYT



# Получение схемы



# Преобразование схемы в StructType

```
{  
  "a": "int32",  
  "b": "string",  
  "c": "uint64",  
  "d": "any",  
  "name.with.dots": "double",  
}
```

—————→ StructType

# Простые типы

```
{  
  "a": "int32",  
  "b": "string",  
  "c": "uint64",  
  "d": "any",  
  "name.with.dots": "double",  
}
```

```
StructType (  
  SF("a", IntegerType),  
  SF("b", StringType),  
  
)
```

# Дополнительные типы

```
{  
  "a": "int32",  
  "b": "string",  
  "c": "uint64",  
  "d": "any",  
  "name.with.dots": "double",  
}
```

```
StructType (  
  SF("a", IntegerType),  
  SF("b", StringType),  
  SF("c", StringType,  
    +original type in meta)  
)
```

# Подсказки для типов

```
{  
  "a": "int32",  
  "b": "string",  
  "c": "uint64",  
  "d": "any",  
  "name.with.dots": "double",  
}
```

+ schema\_hint

"d": ArrayType(StringType)

```
StructType (  
  SF("a", IntegerType),  
  SF("b", StringType),  
  SF("c", StringType, +meta)  
  SF("d", ArrayType(StringType)),  
)
```

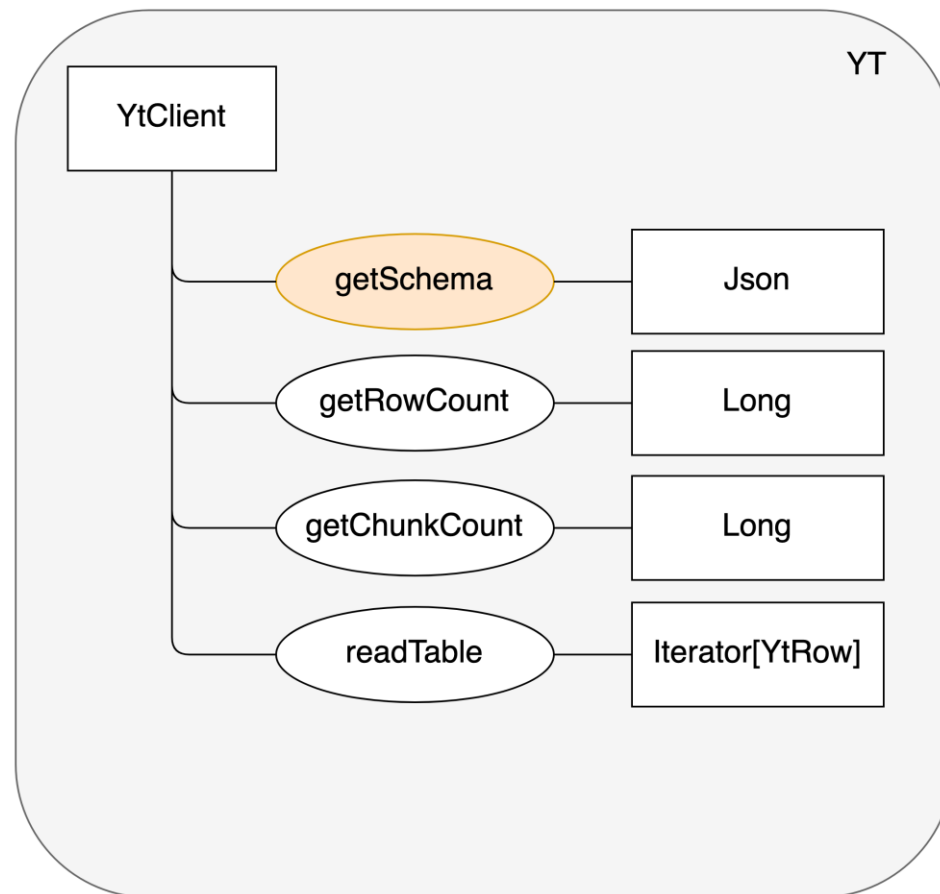
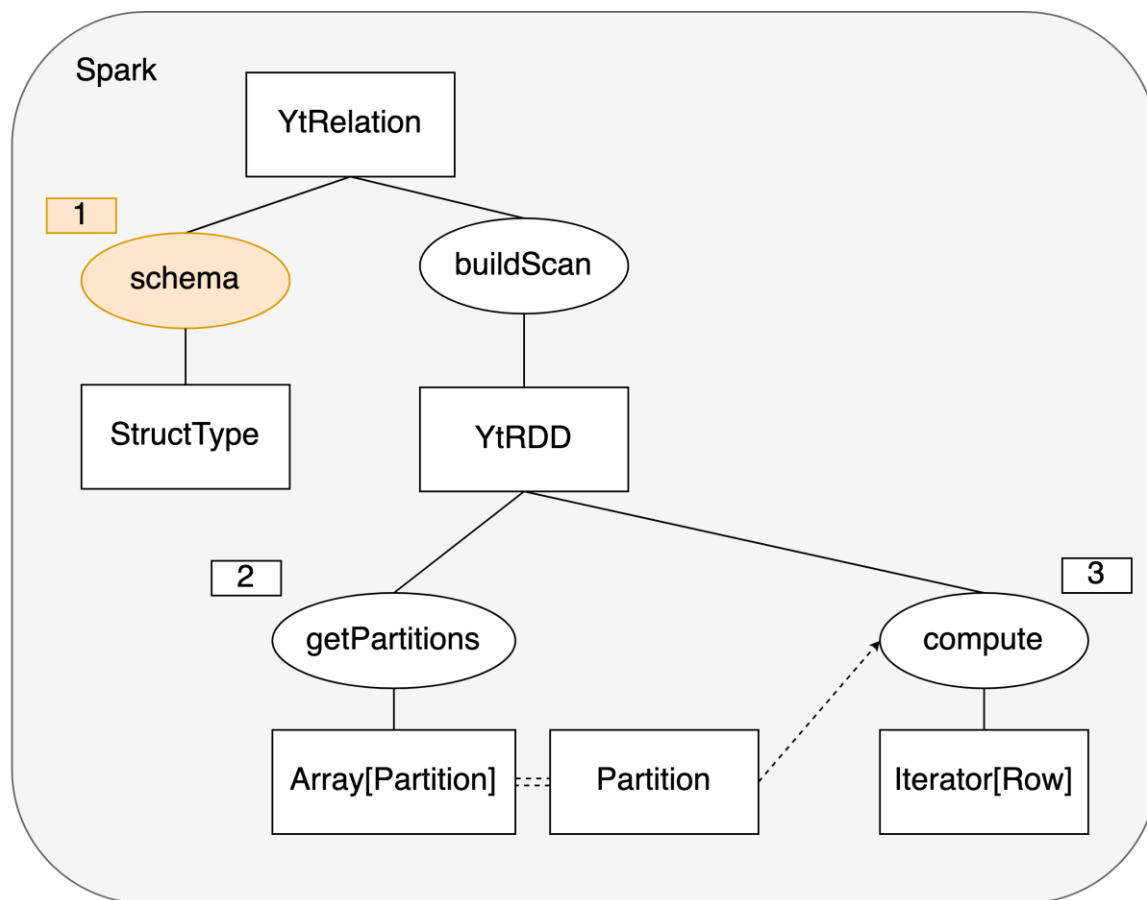


# Переименование

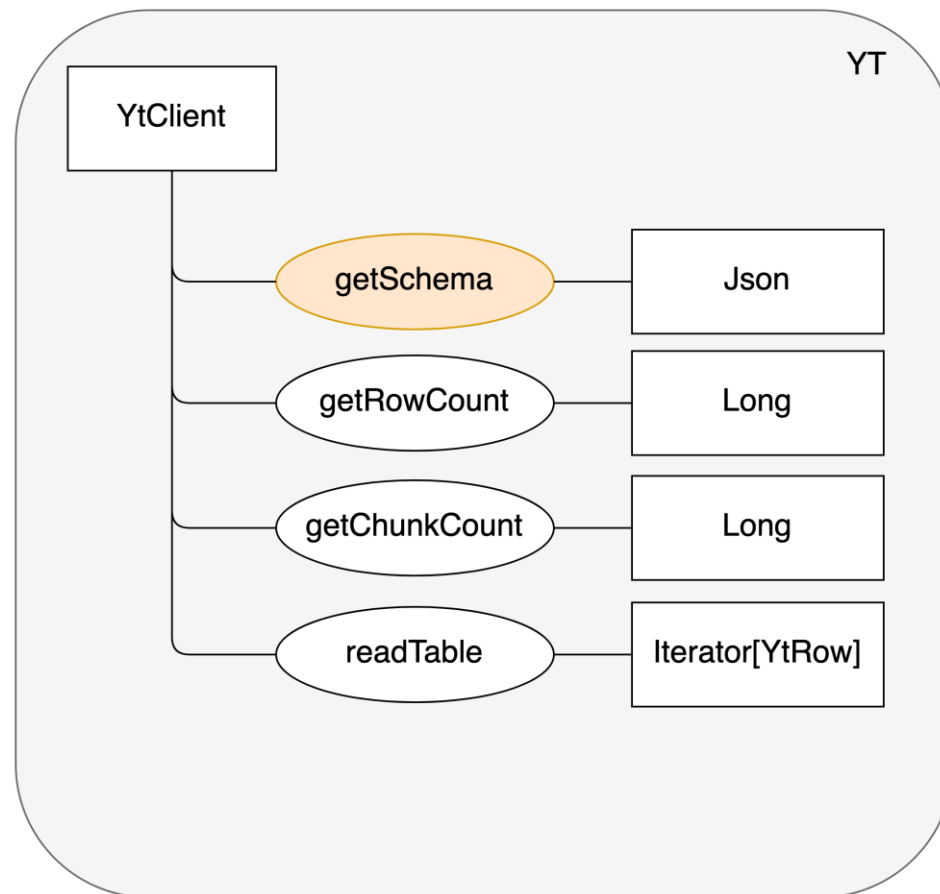
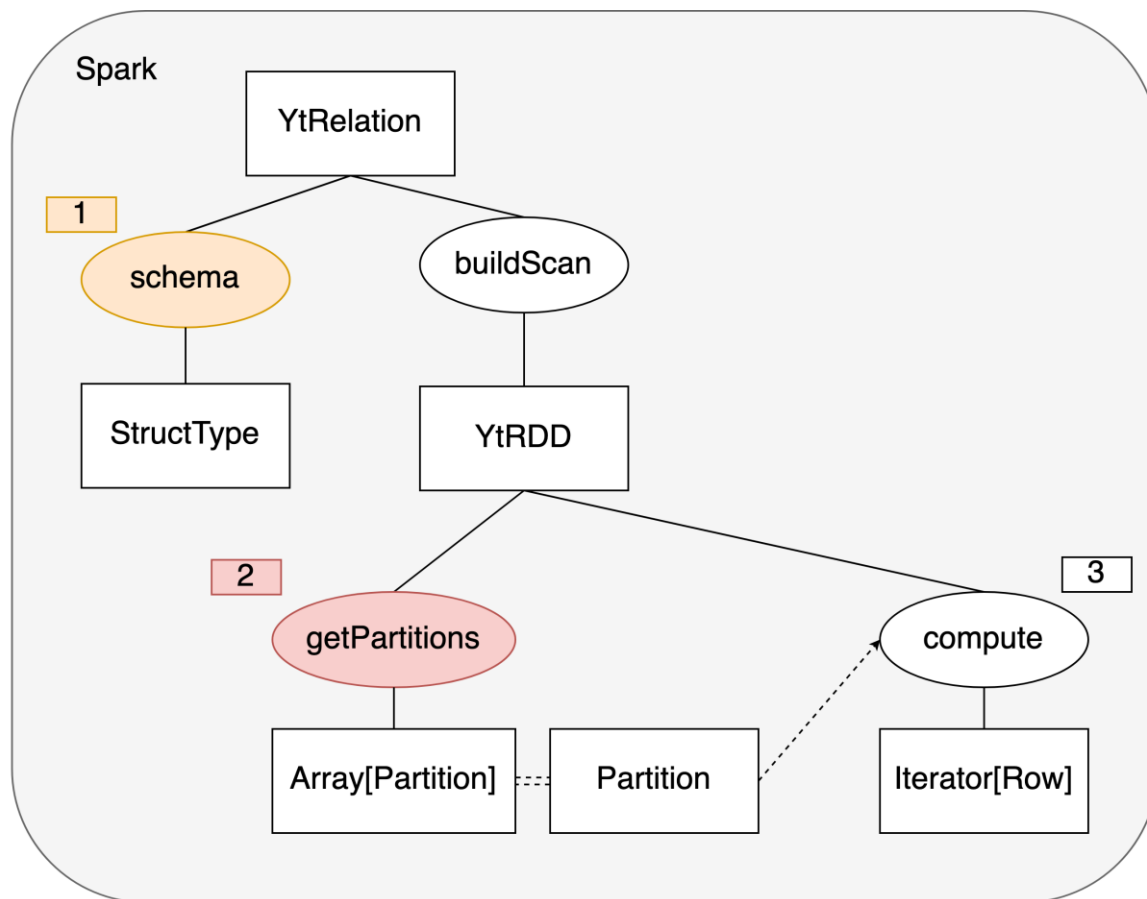
```
{  
  "a": "int32",  
  "b": "string",  
  "c": "uint64",  
  "d": "any",  
  "name.with.dots": "double",  
}
```

```
StructType (  
  SF("a", IntegerType),  
  SF("b", StringType),  
  SF("c", StringType, +meta)  
  SF("d", ArrayType(StringType)),  
  SF("name_with_dots",  
    DoubleType,  
    + original name in meta)  
)
```

# Получение схемы



# Получение партиций



# Таблица в YT

## User level

#	a	b	c
0	"cat"	true	1
1	"dog"	true	20
2	"otter"	false	350

### Attributes:

schema  
a - String  
b - Boolean  
c - Integer

rowCount

## Storage level

### Columnar chunk 1

rows 0 - 200

### Columnar chunk 2

rows 201 - 400

### Columnar chunk 3

rows 401 - 600

### Attributes:

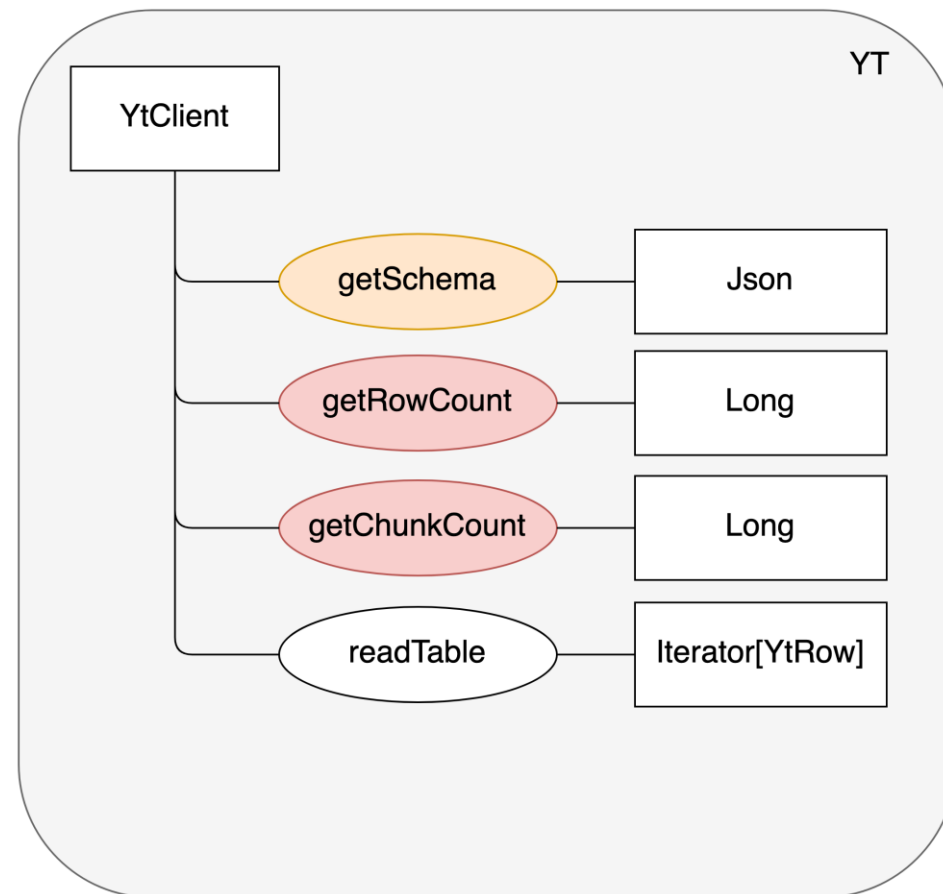
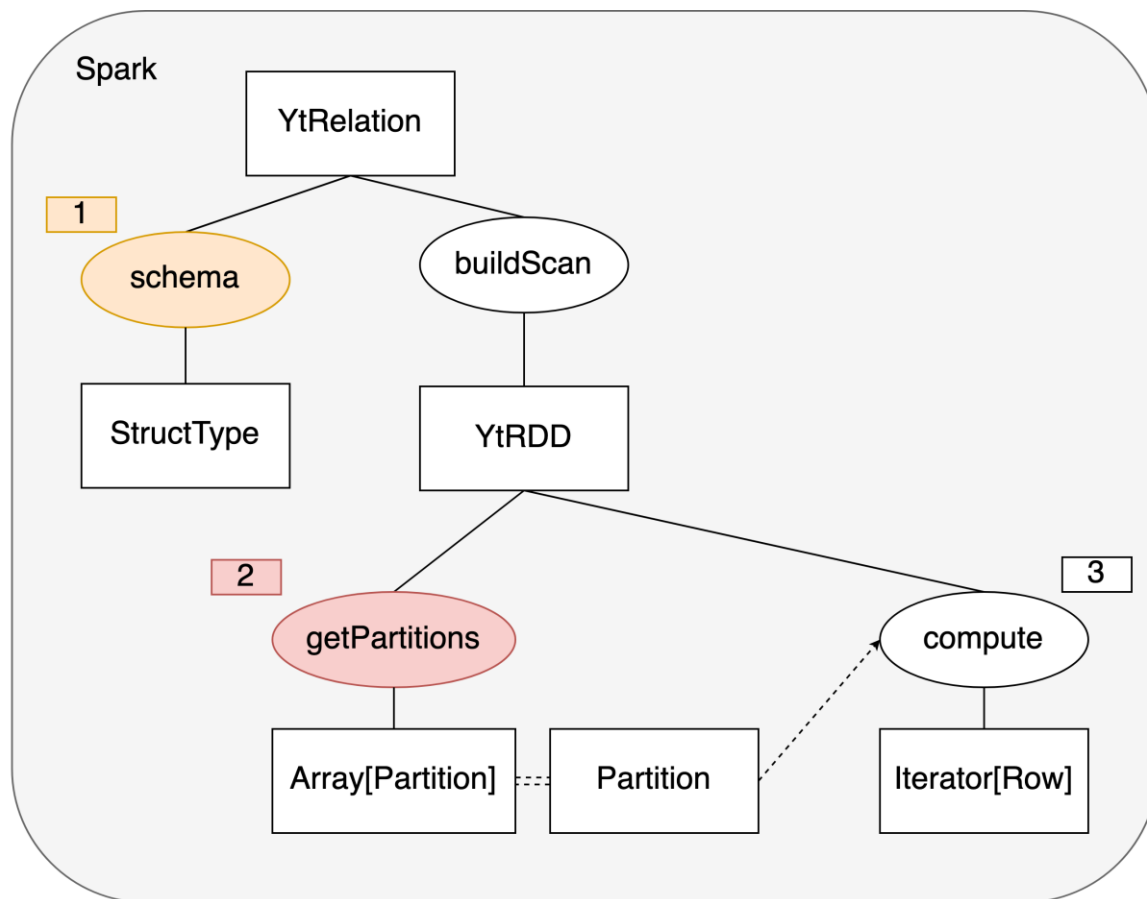
schema  
a - String  
b - Boolean  
c - Integer

rowCount

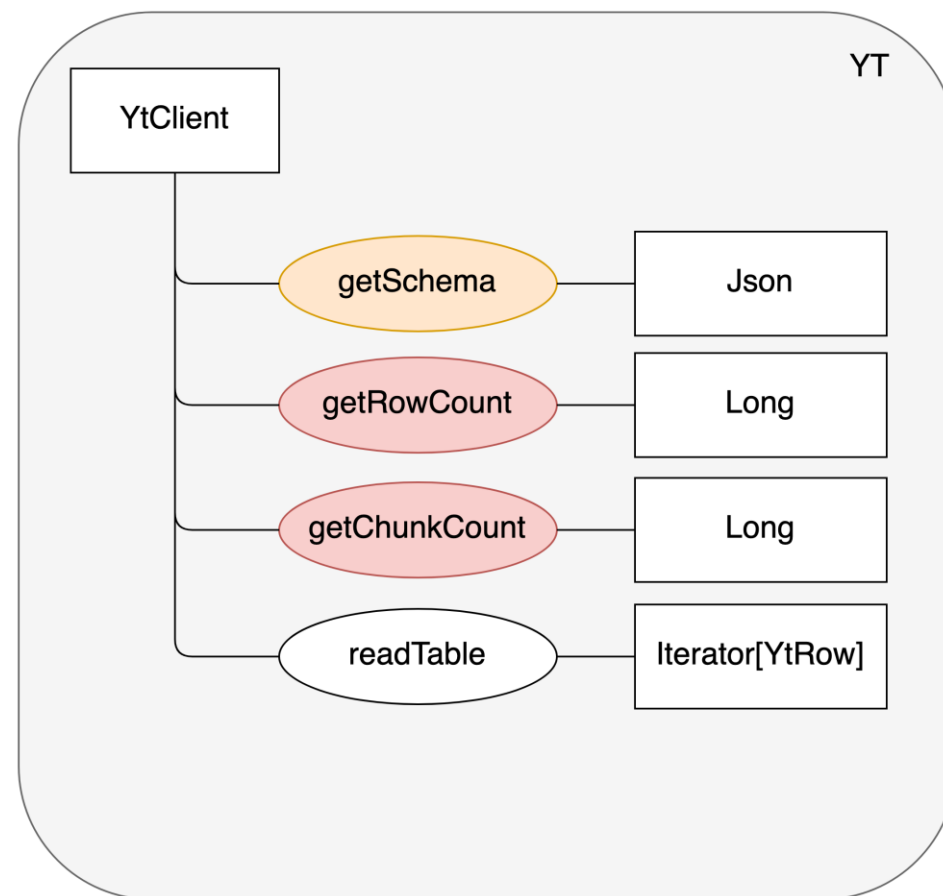
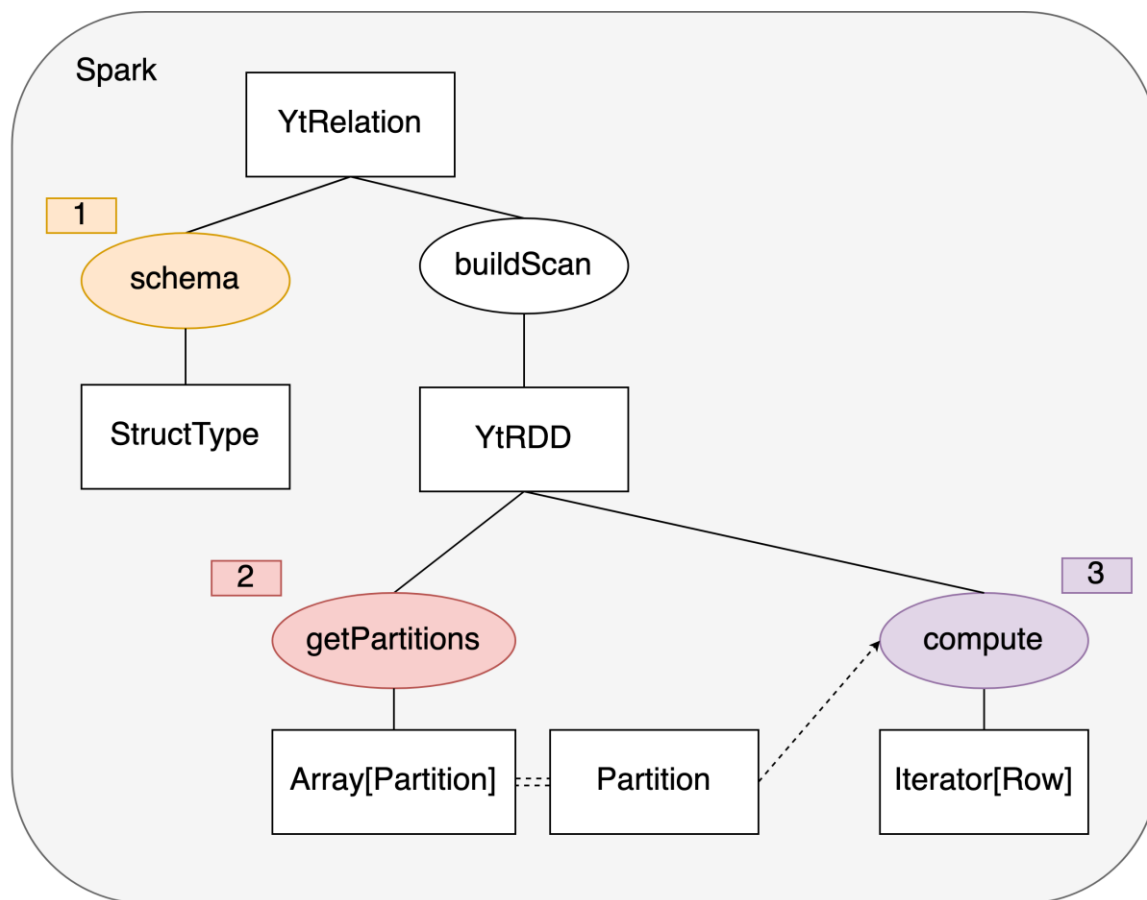
# Разбиение таблицы на партиции

- Можем узнать точное количество строк в таблице из метаданных
- Можем запрашивать диапазоны строк при чтении
- Бывают проблемы с равномерностью распределения

# Получение партиций



# Чтение партиции

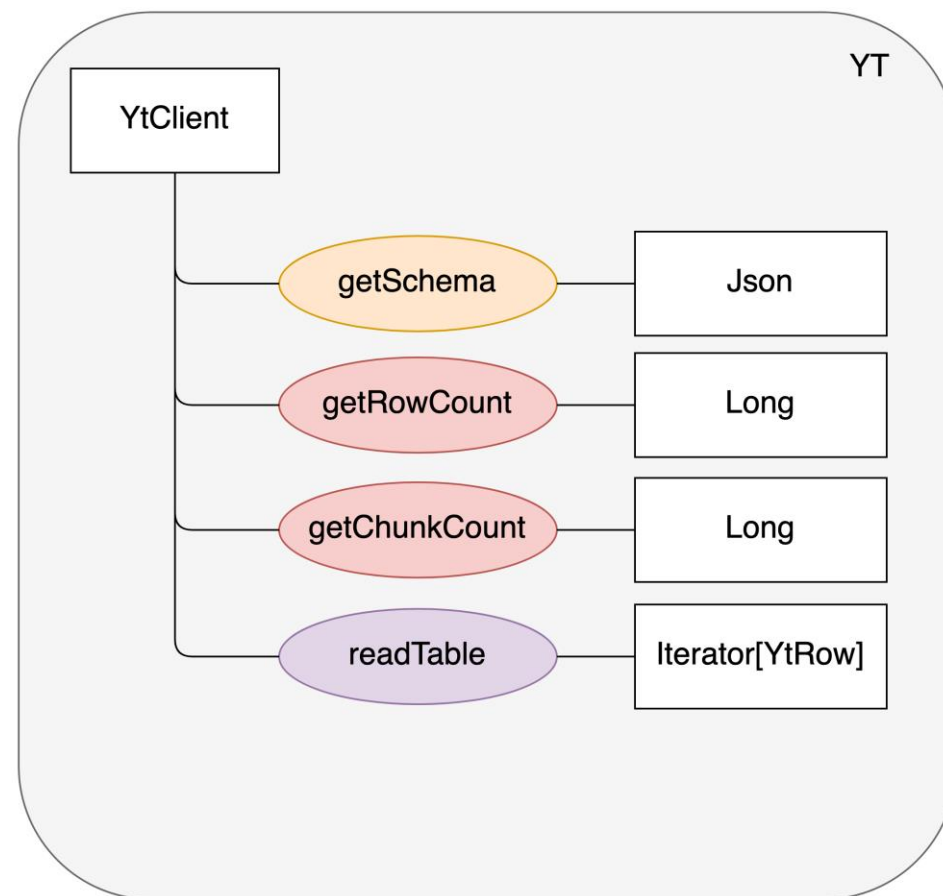
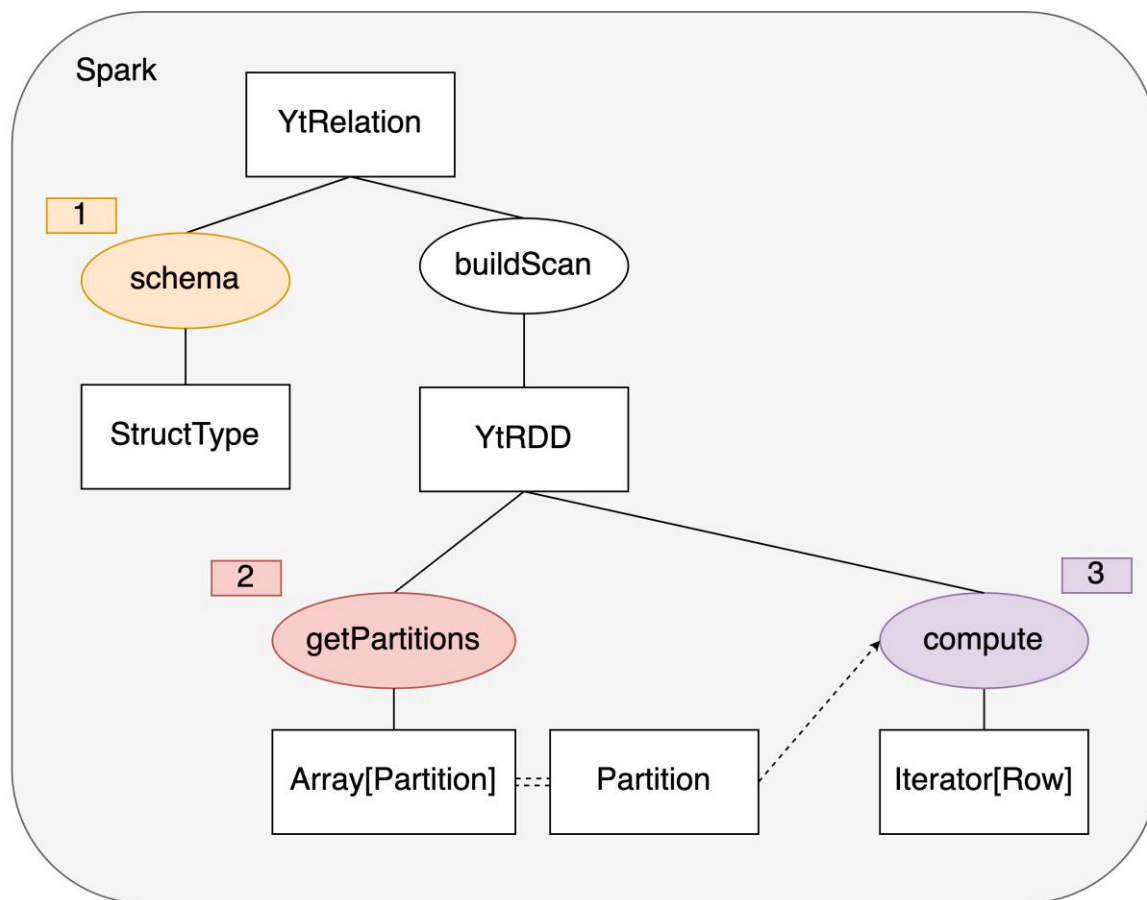


# Чтение партиции из YT

- API отдаёт строки таблицы в своём формате YtRow
- Нужно преобразовывать строки в спарковый Row
- В этот момент нужно сделать преобразование значений под спарковую систему типов

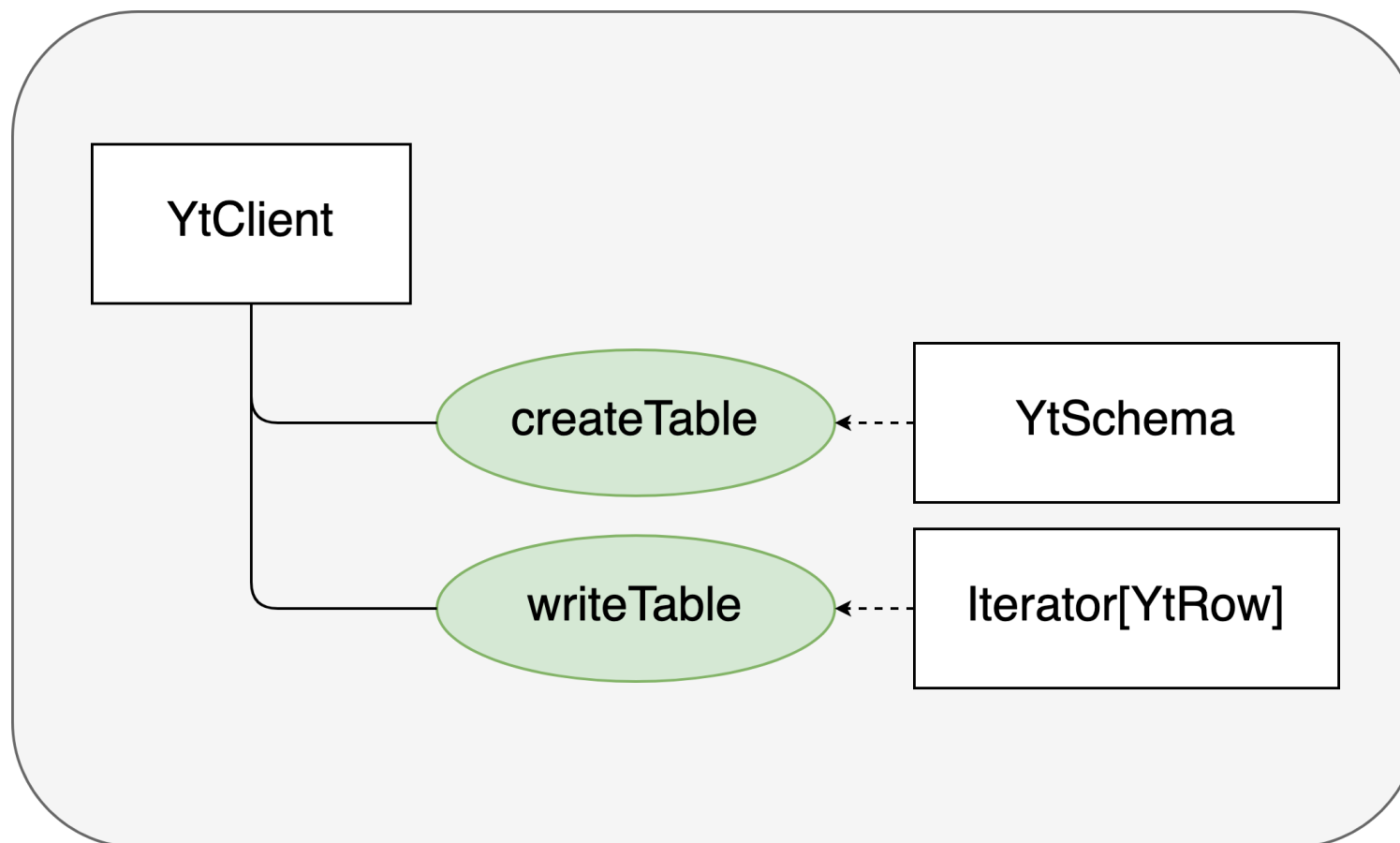


# Чтение

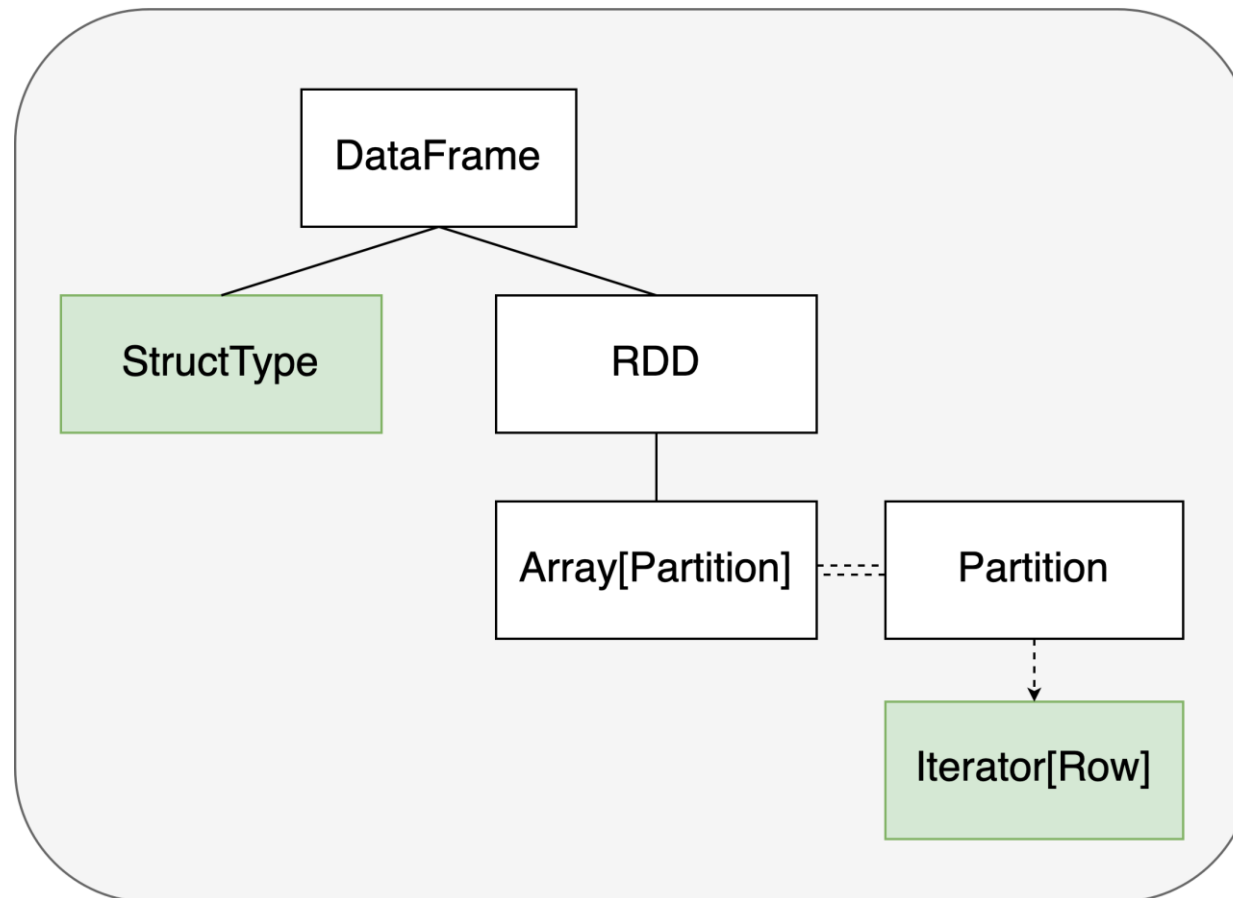


Простое решение: запись

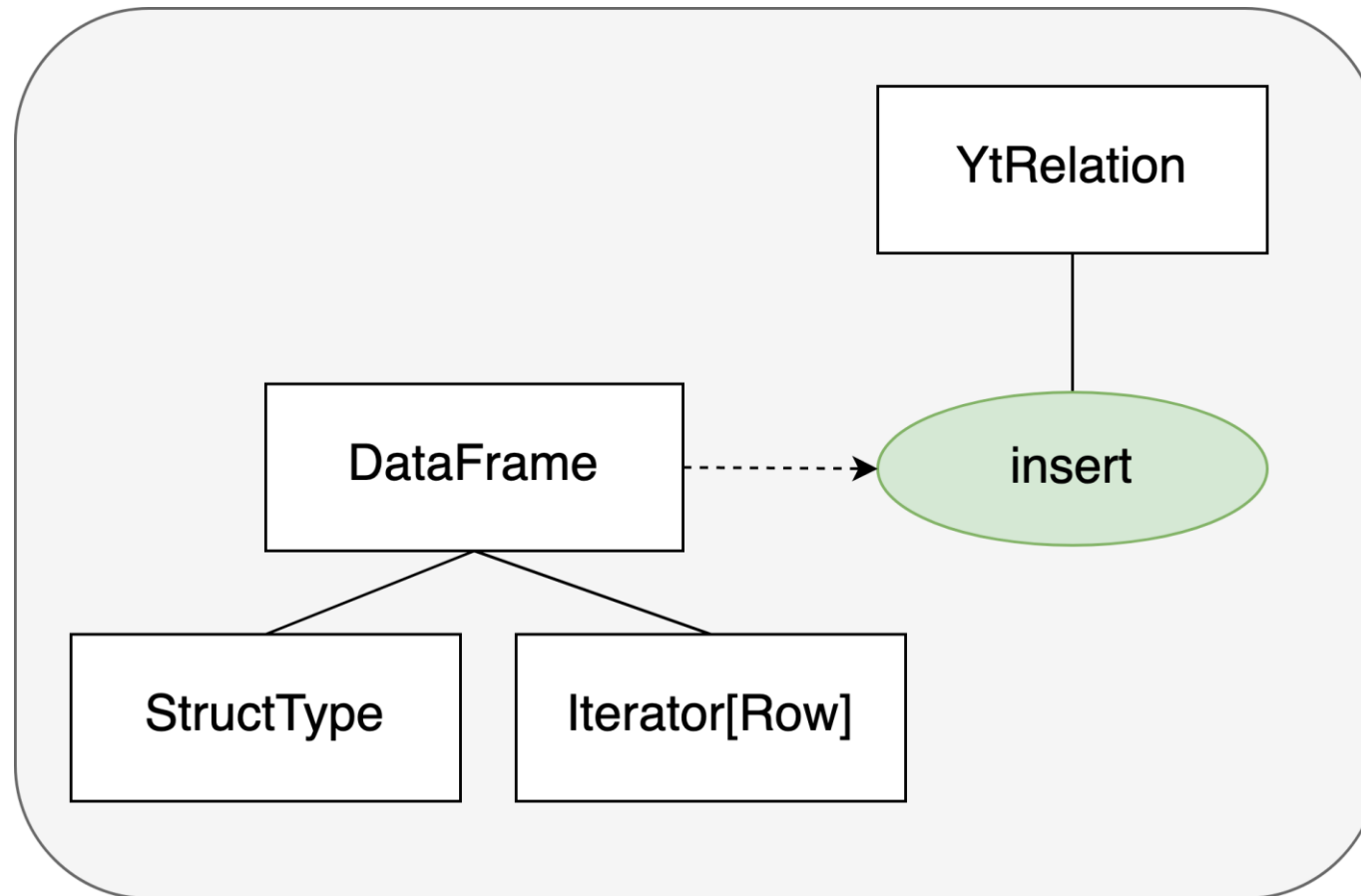
# YT API



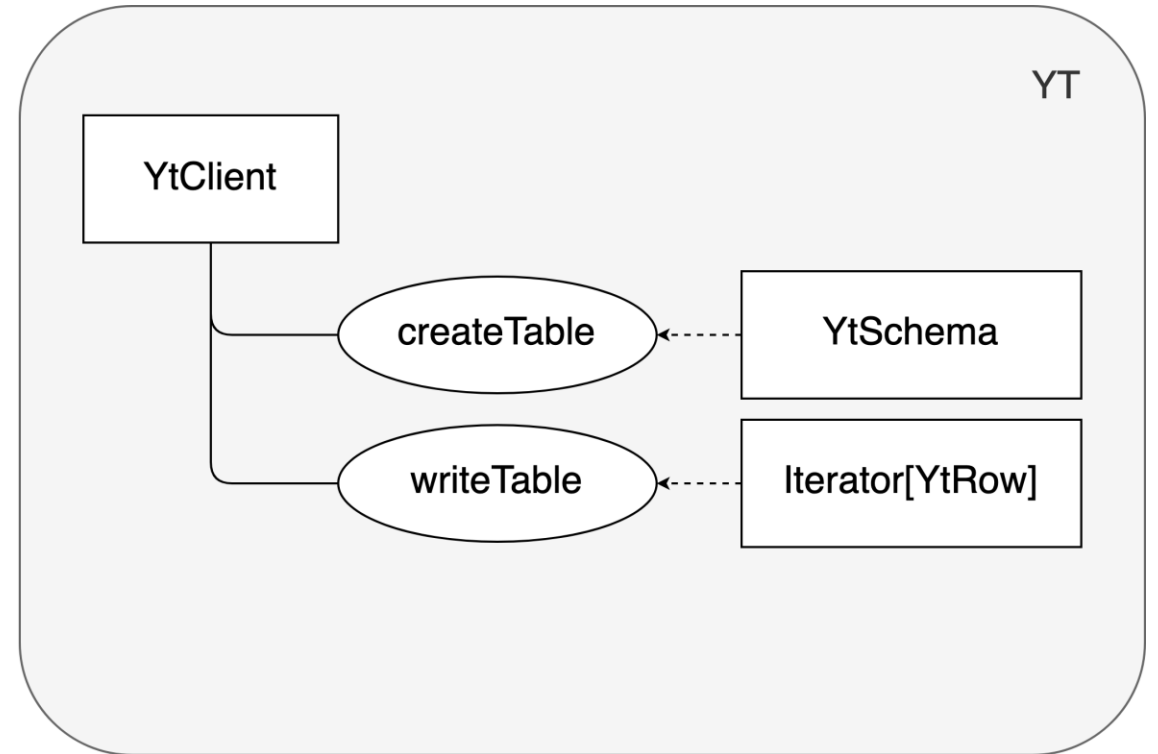
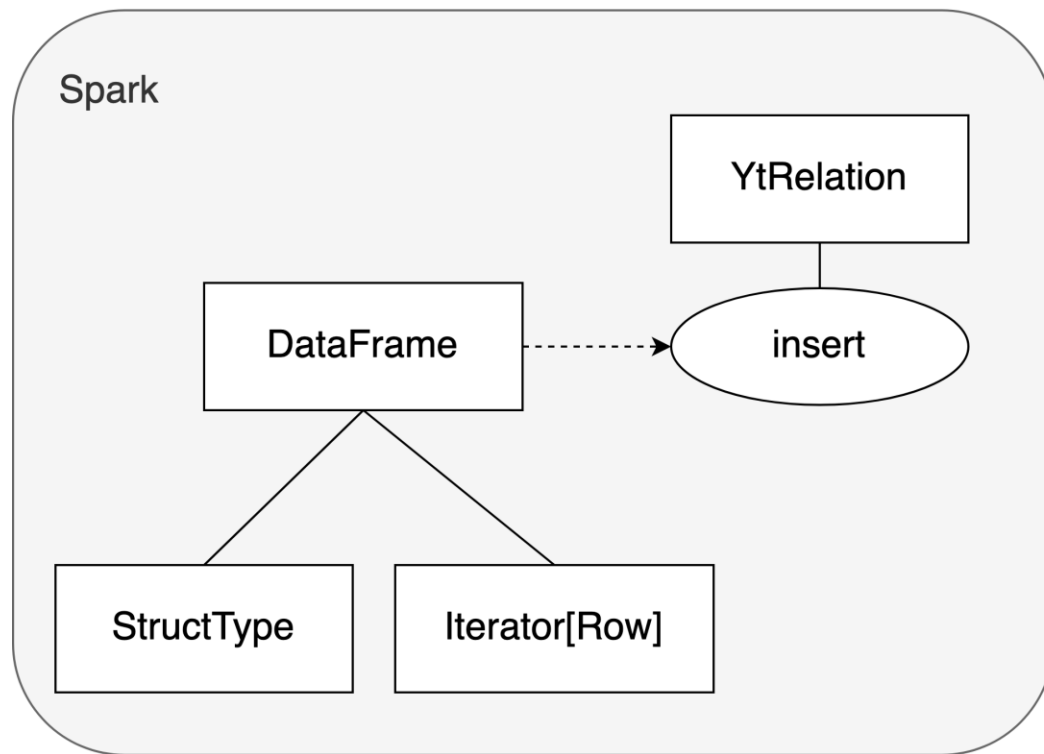
# Spark SPI



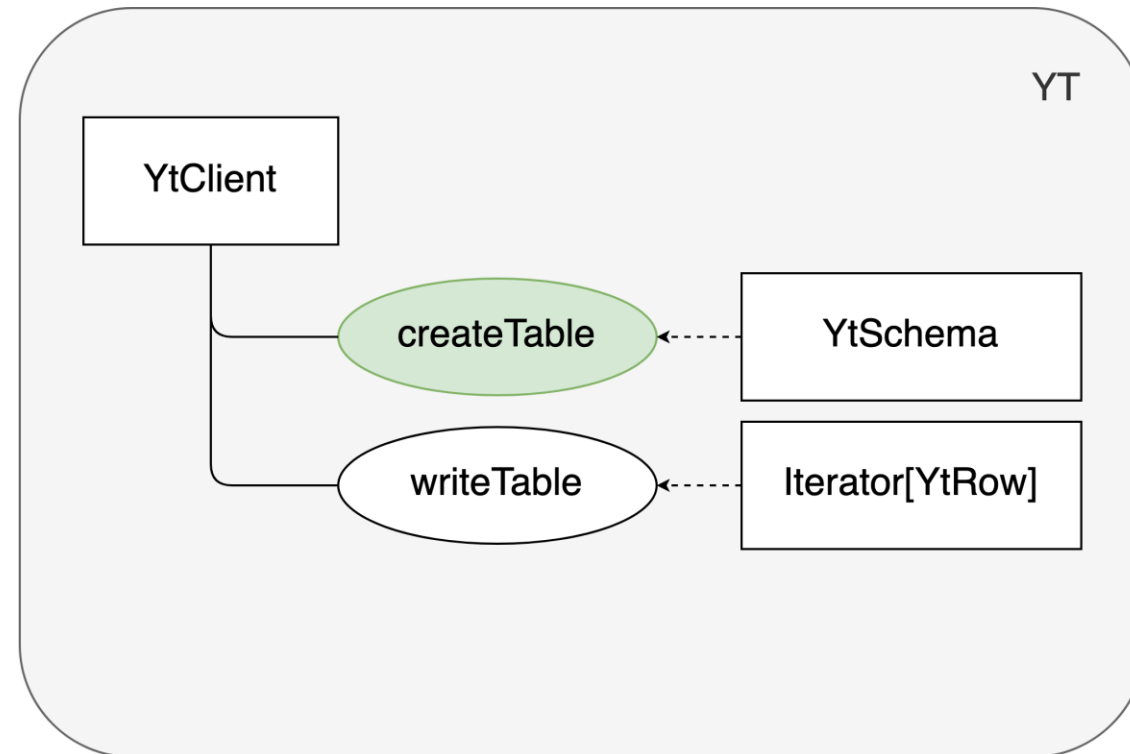
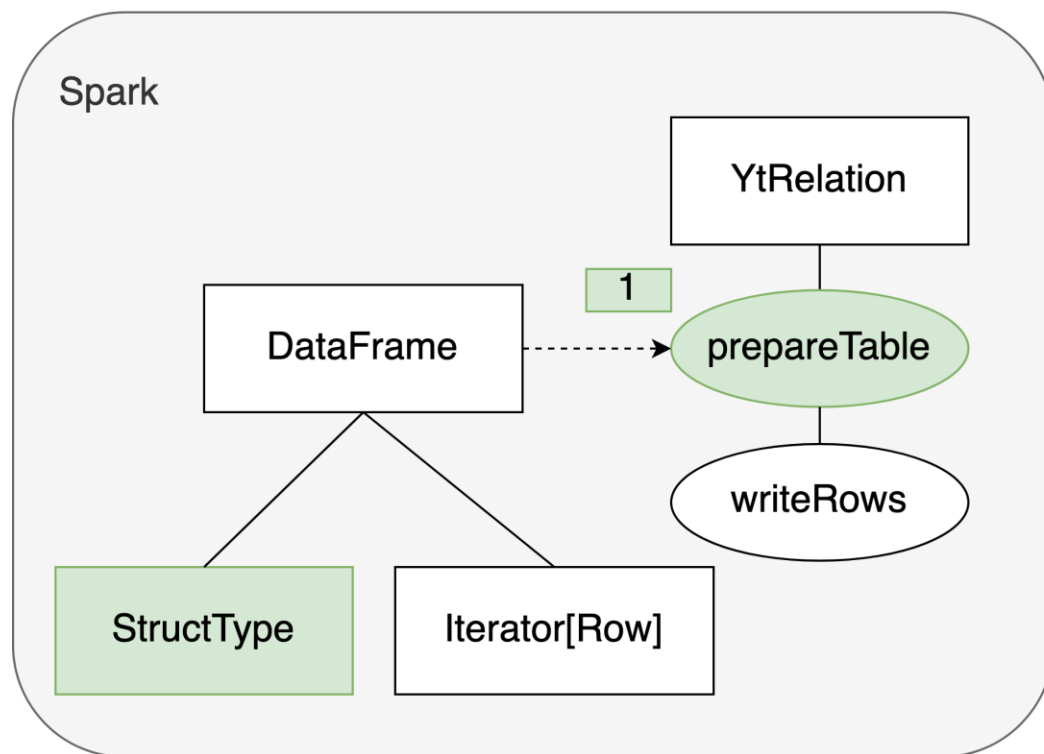
# Spark SPI



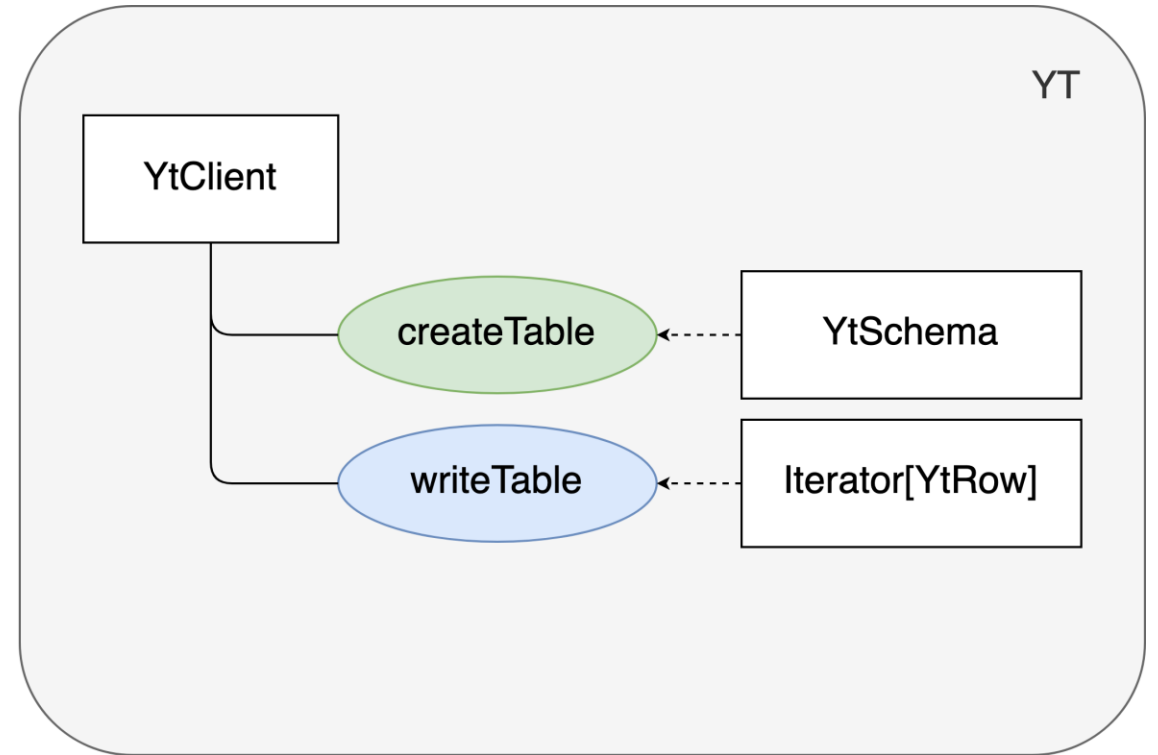
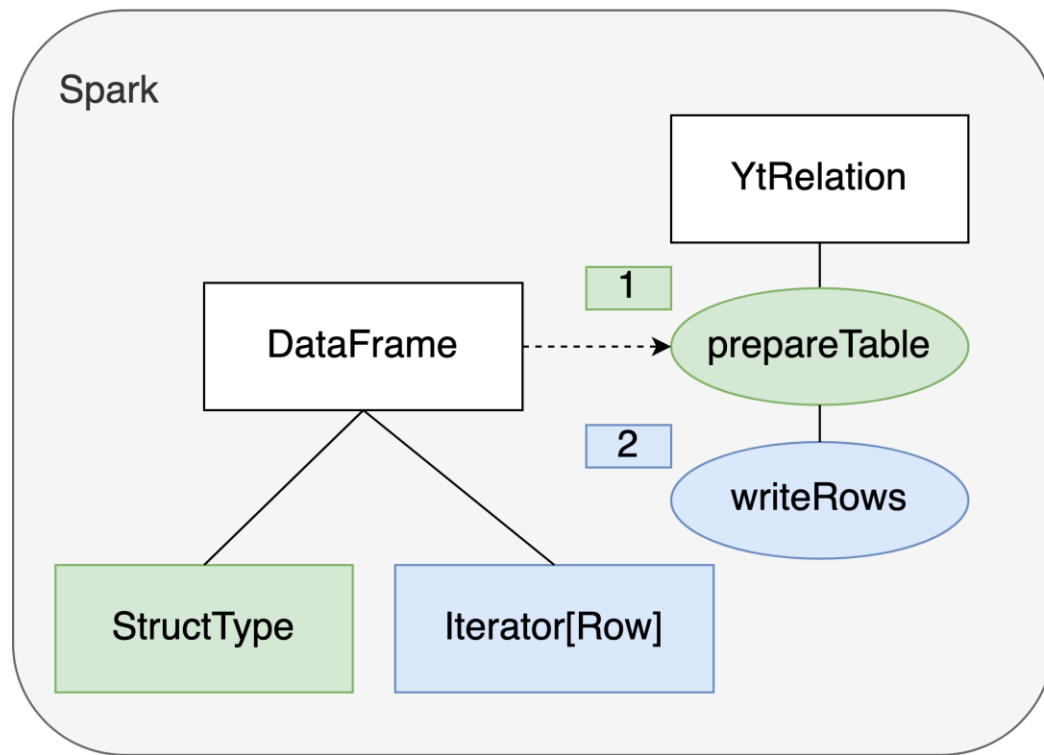
# SPYT



# Подготовка записи

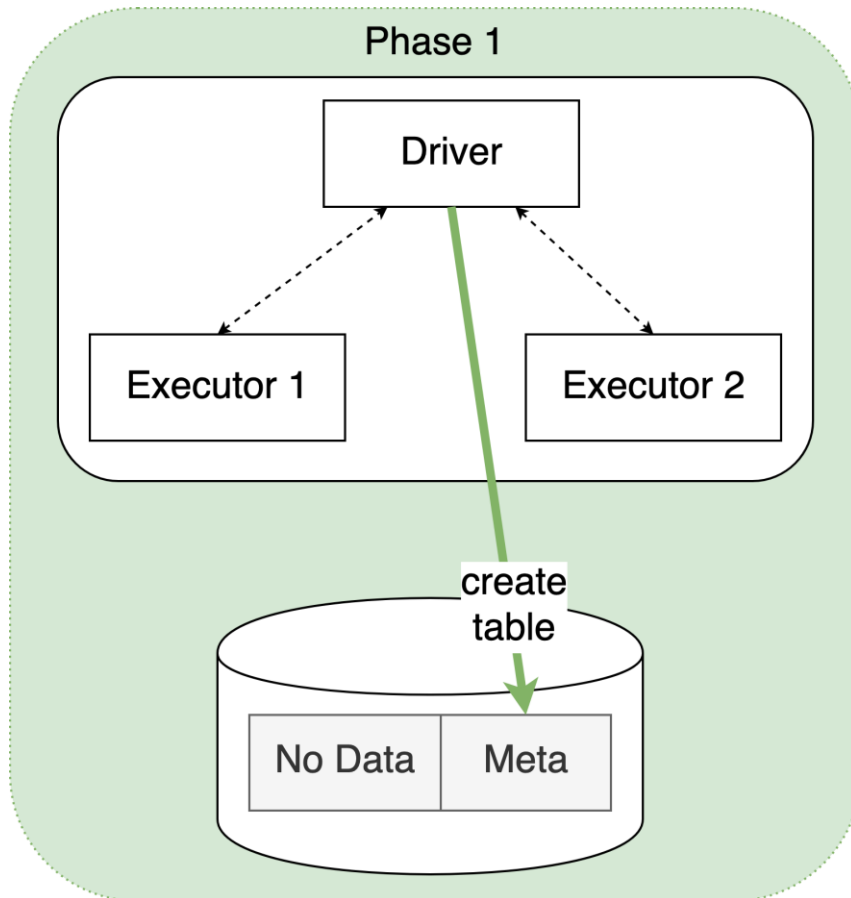


# Запись партиций

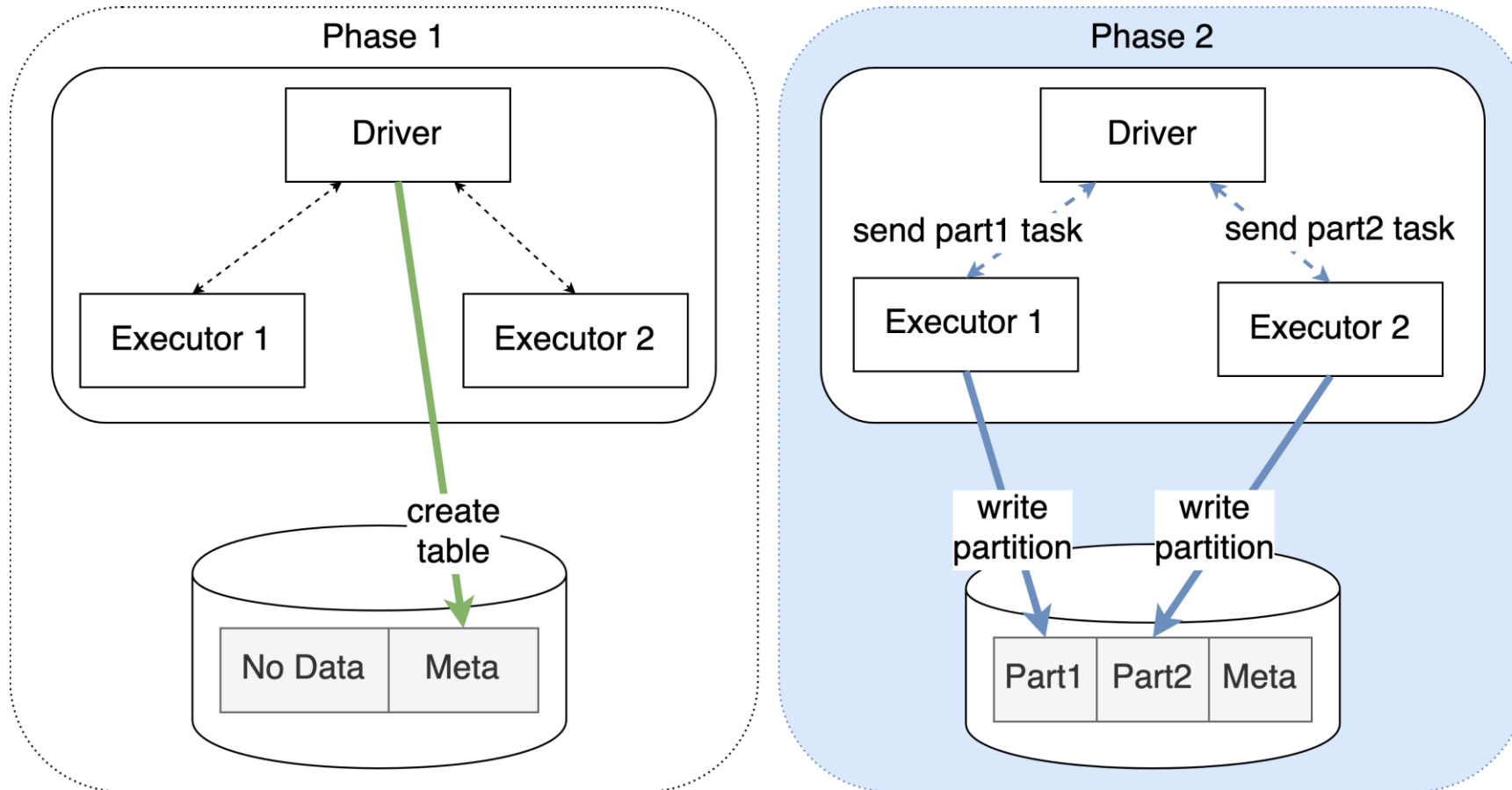




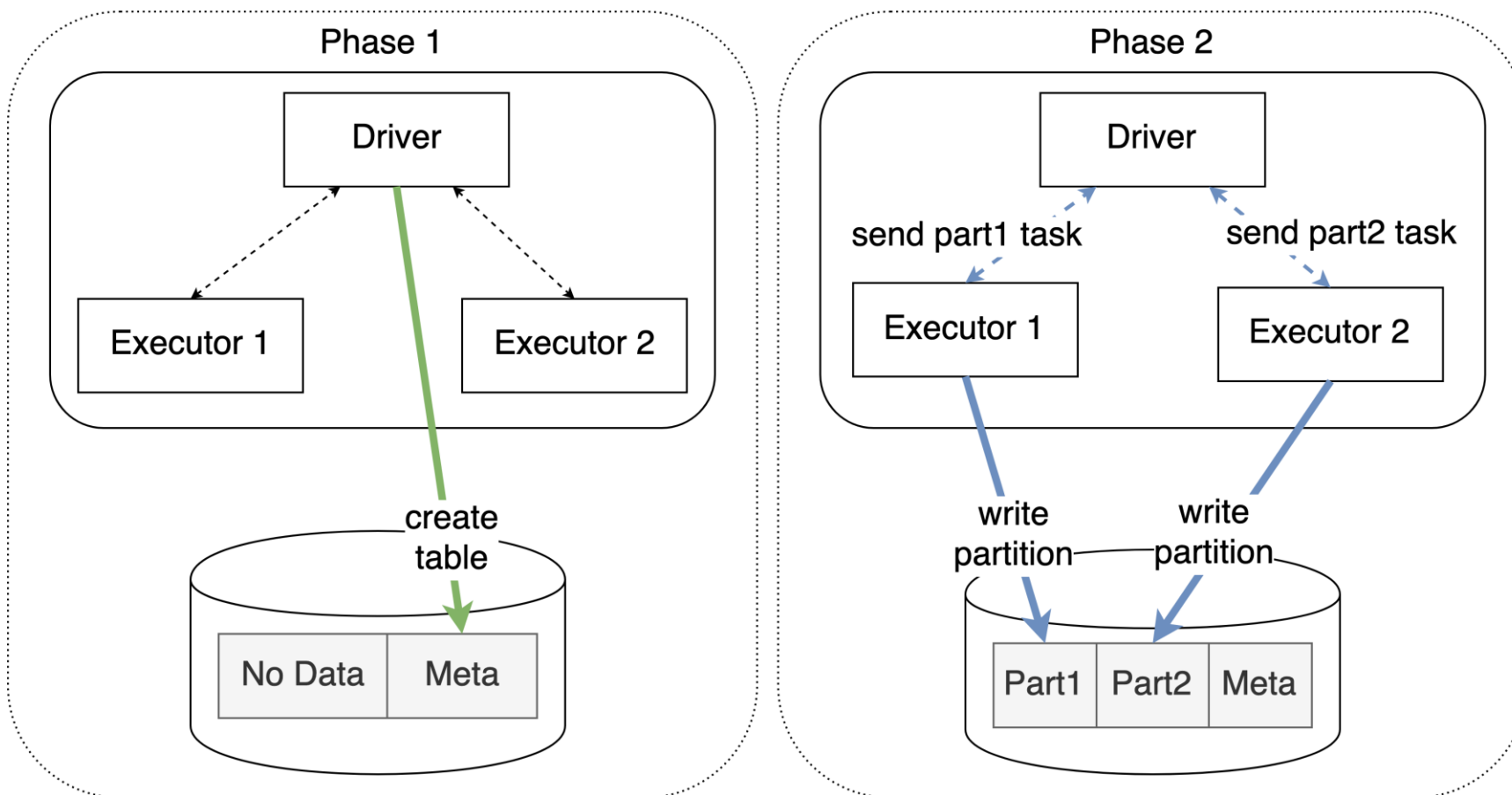
# Подготовка таблицы



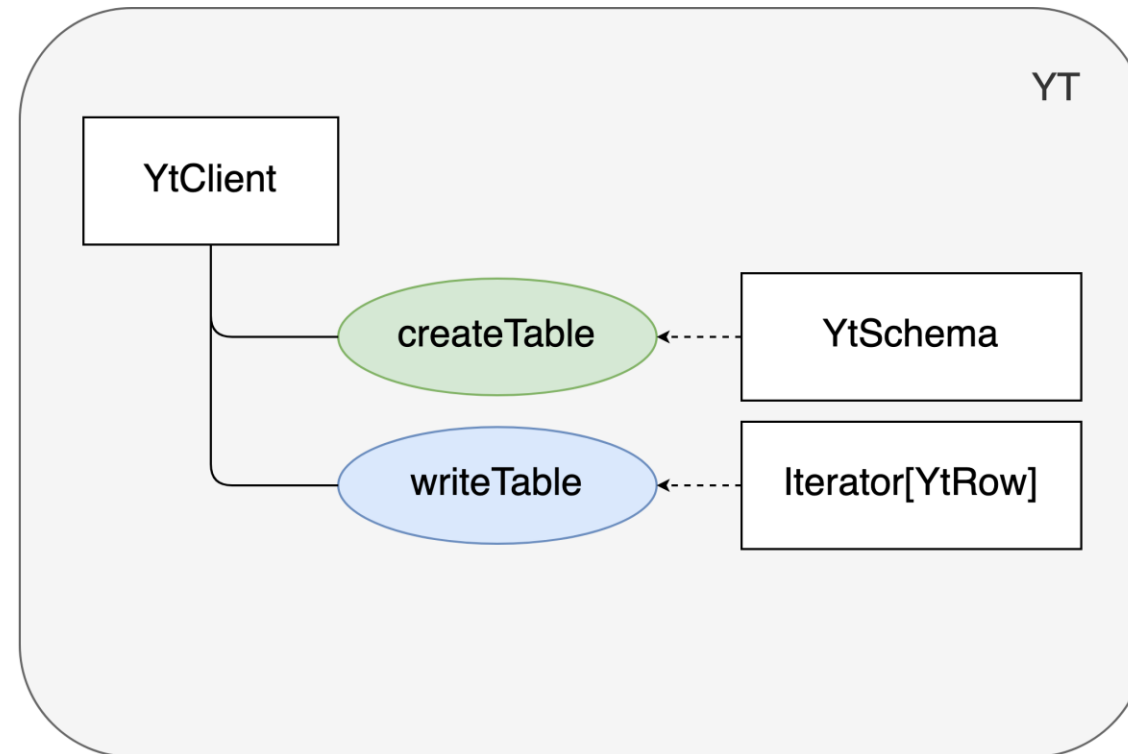
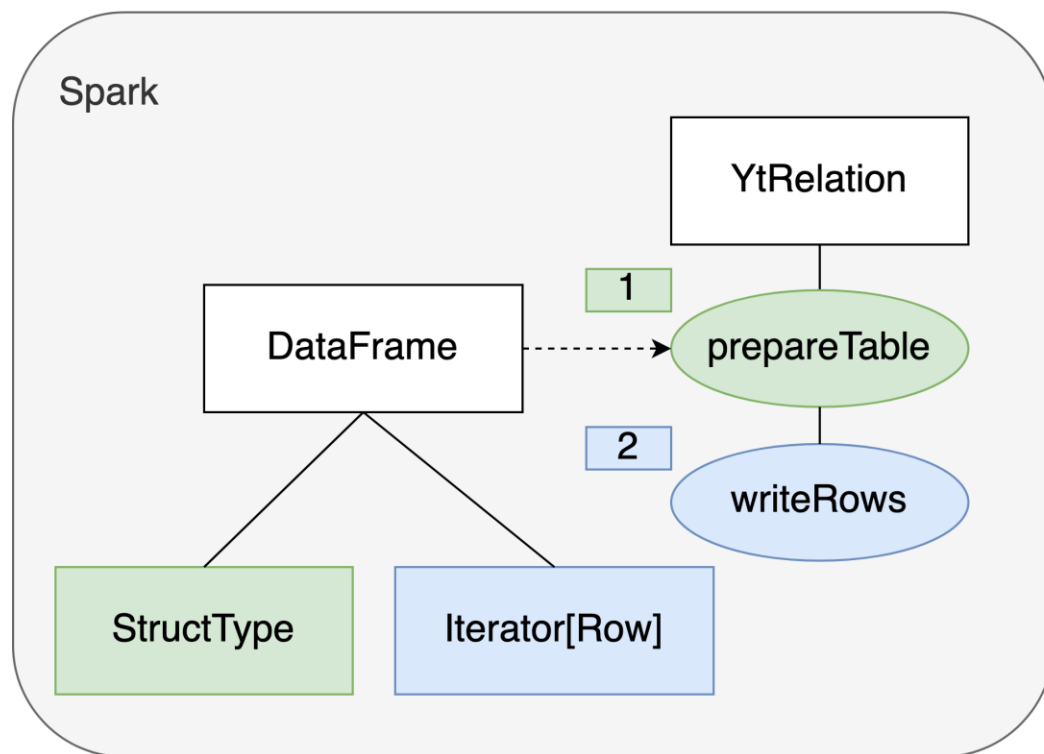
# Запись партиций



# Как это выполняется



# Как это выглядит в коде



# Особенности записи

- Обратное преобразование схемы
- Обратное преобразование строчек таблицы
- Метод `writeRows` будет вызываться параллельно с разных экзекуторов

# Итог простого решения

- Чтение:
  - Запрос схемы с драйвера
  - Разделение на партиции на драйвере
  - Чтение партиции на экзекьюторе
- Запись:
  - Создание таблицы на драйвере
  - Запись партиций с экзекьюторов параллельно

Подробное выступление Jacek Laskowski: <https://www.youtube.com/watch?v=vfd83ELIMfc>

# Что не так с простым подходом?

# Что не так с простым подходом

- Нет чтения батчами



# Таблица в YT

## User level

#	a	b	c
0	"cat"	true	1
1	"dog"	true	20
2	"otter"	false	350

### Attributes:

schema  
a - String  
b - Boolean  
c - Integer  
  
rowCount

## Storage level

### Columnar chunk 1

rows 0 - 200

### Columnar chunk 2

rows 201 - 400

### Columnar chunk 3

rows 401 - 600

### Attributes:

schema  
a - String  
b - Boolean  
c - Integer  
  
rowCount

# Батчи в УТ

Storage level

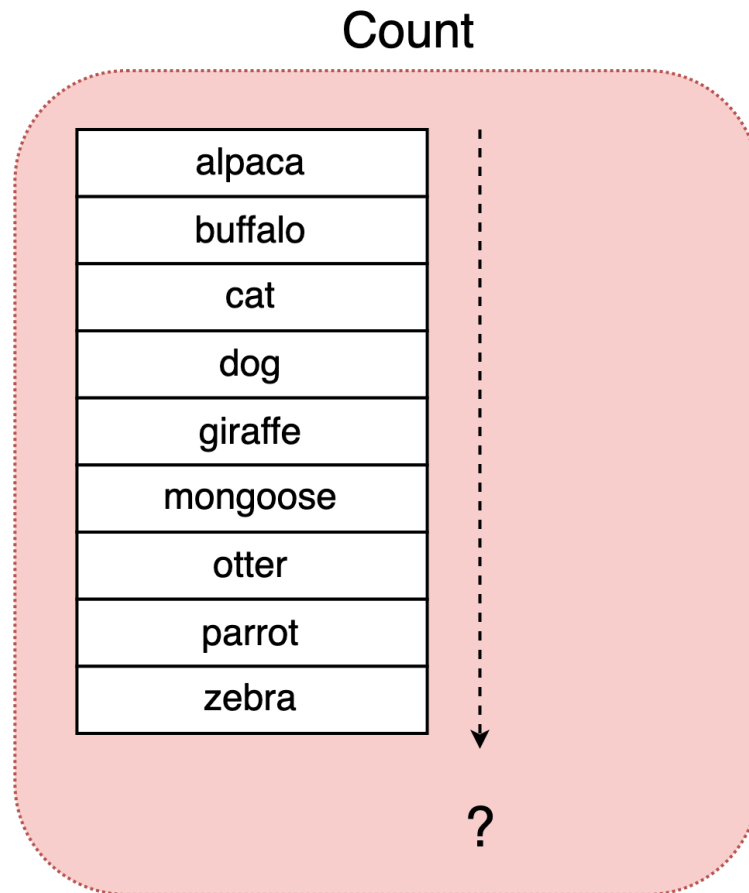
Columnar chunk 1 rows 0 - 200	Chunk meta  200 rows a - ["alpaca", "cat"]
Columnar chunk 2 rows 201 - 400	Chunk meta  200 rows a - ["dog", "mongoose"]
Columnar chunk 3 rows 401 - 600	Chunk meta  200 rows a - ["otter", "zebra"]

## Attributes:

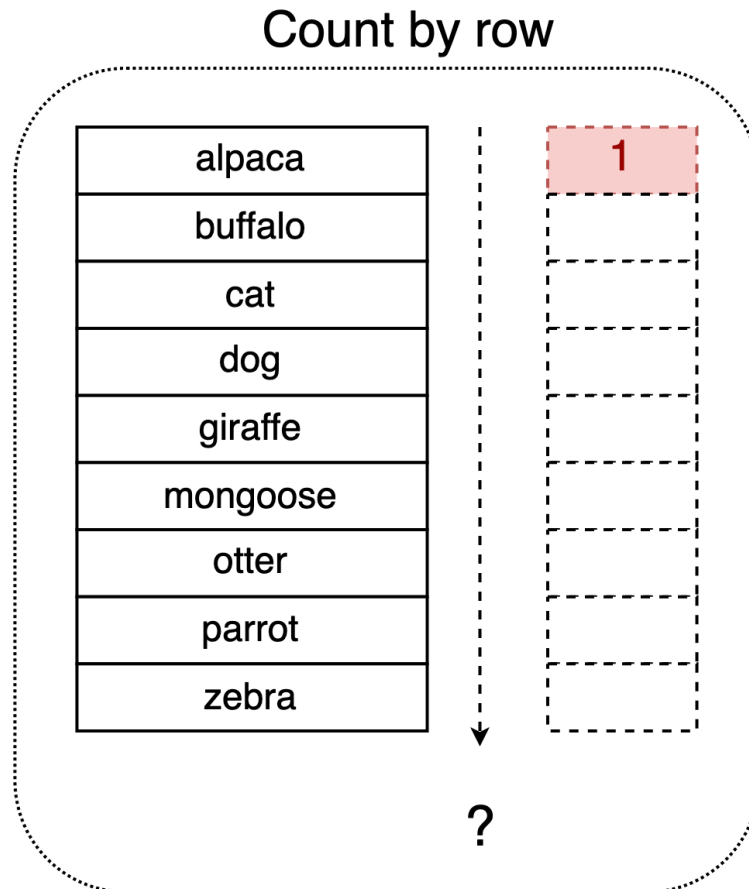
schema  
a - String  
b - Boolean  
c - Integer  
  
rowCount

...

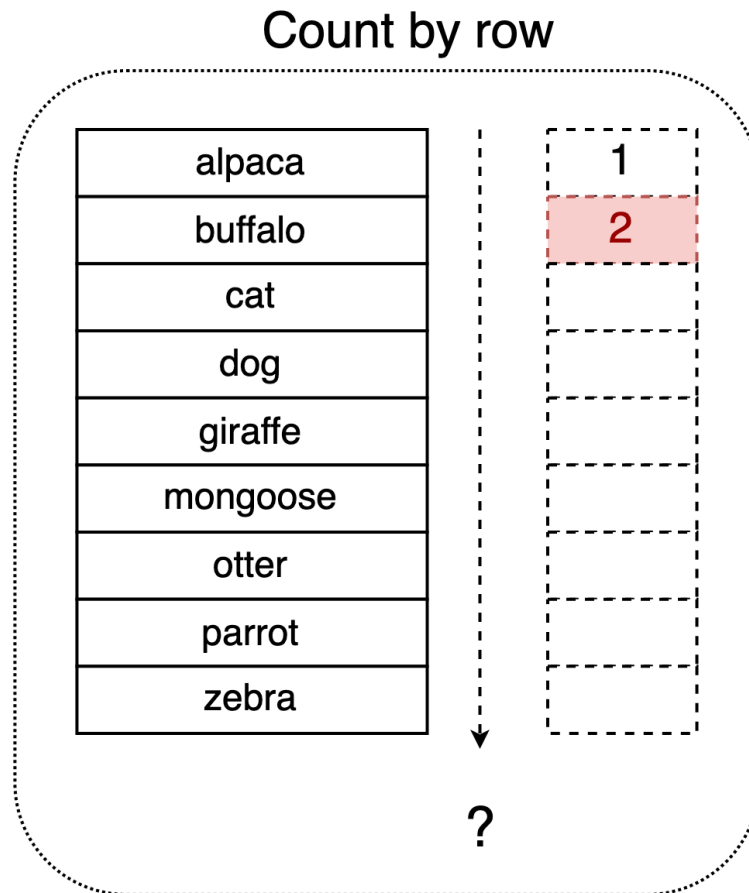
# Батчи в Spark



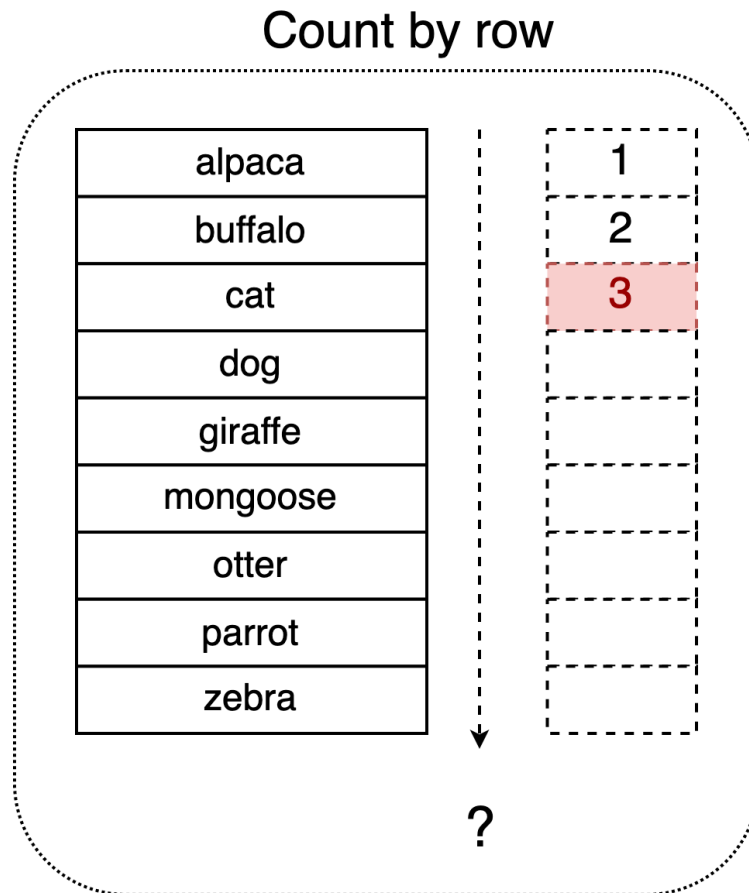
# Построчный count



# Построчный count

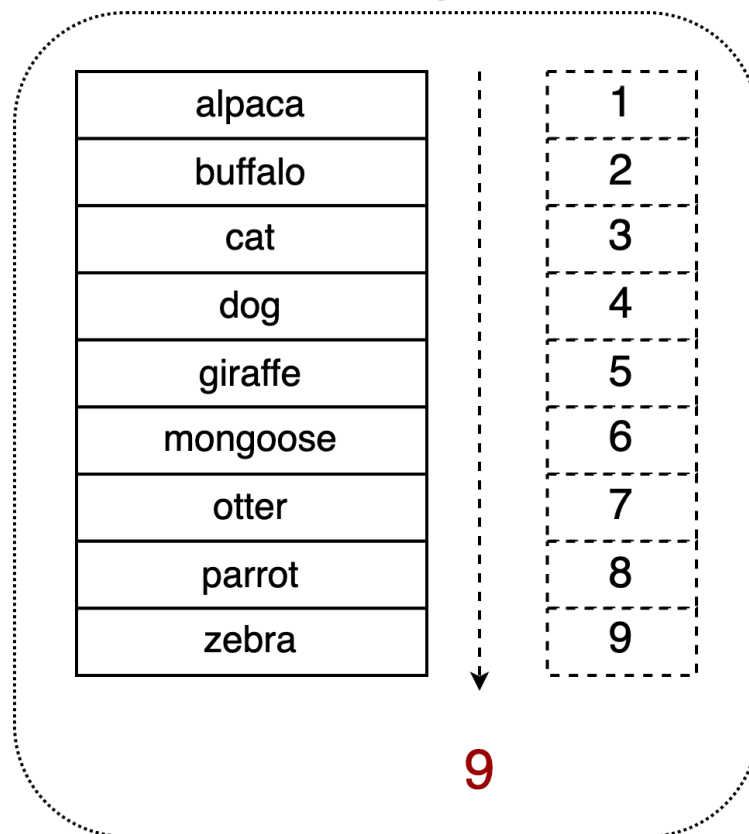


# Построчный count



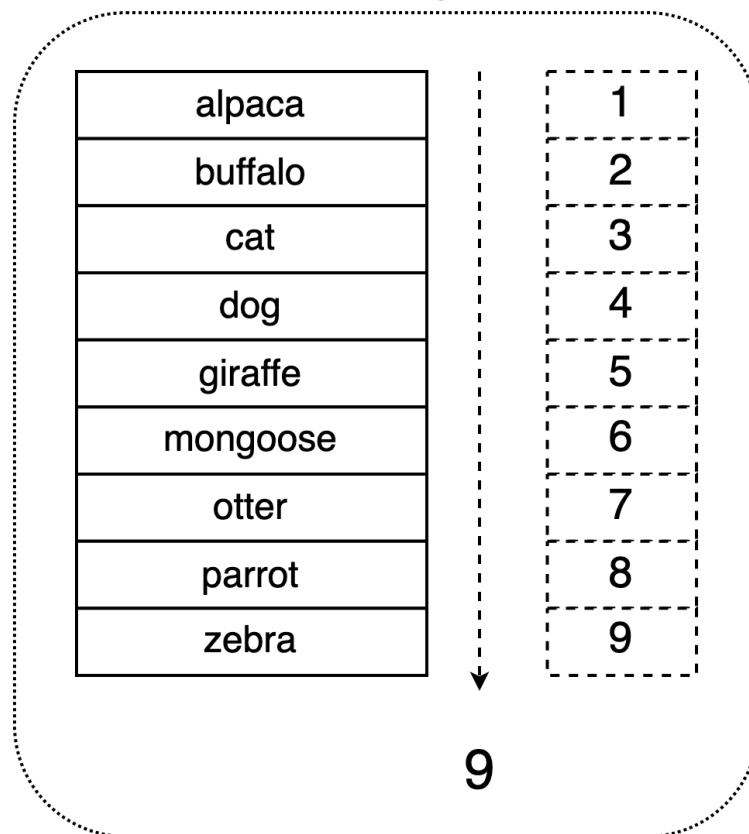
# Построчный count

Count by row

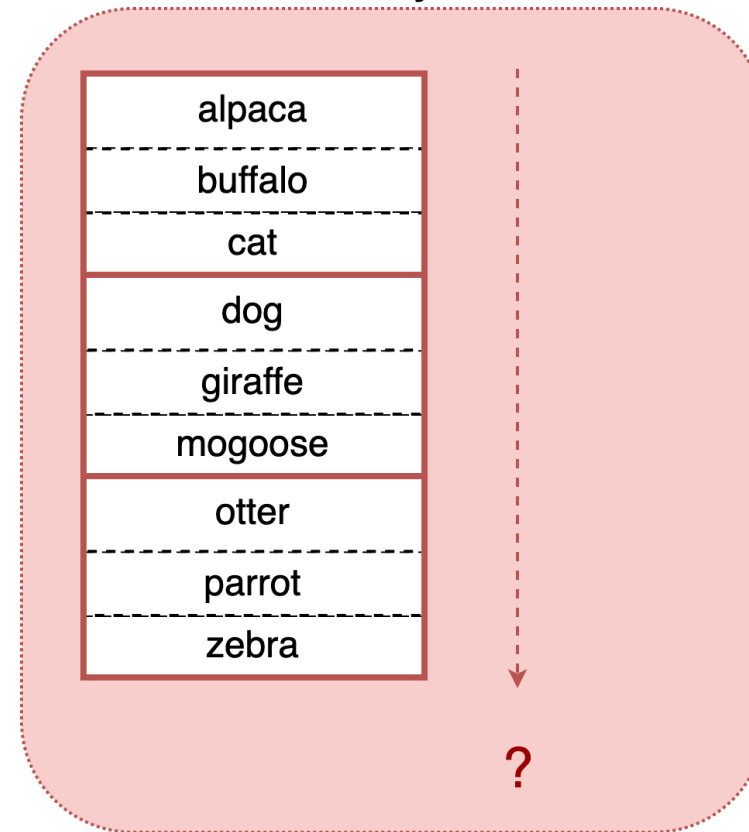


# Побатчевый count

Count by row



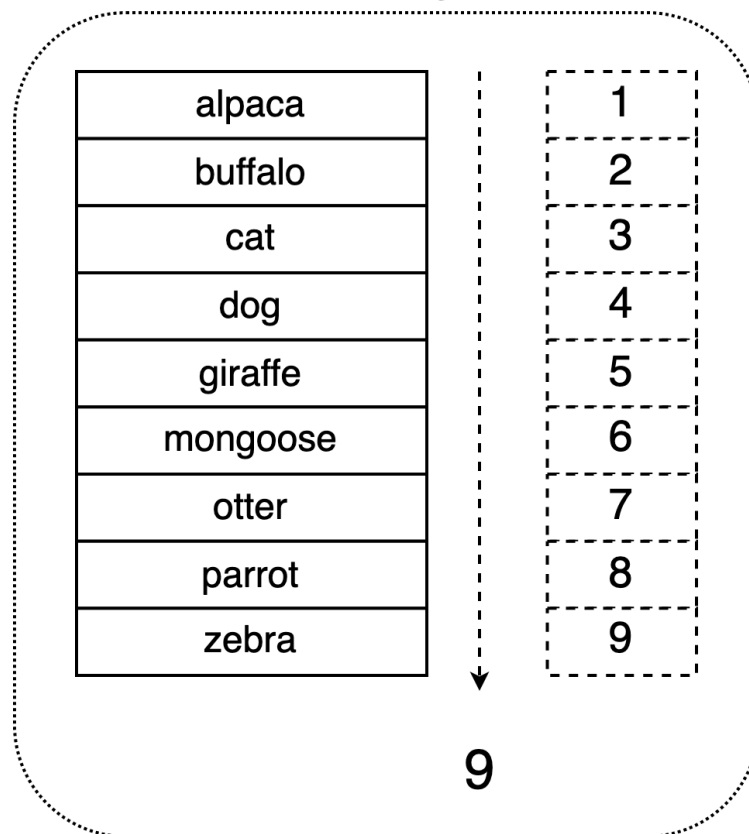
Count by batch



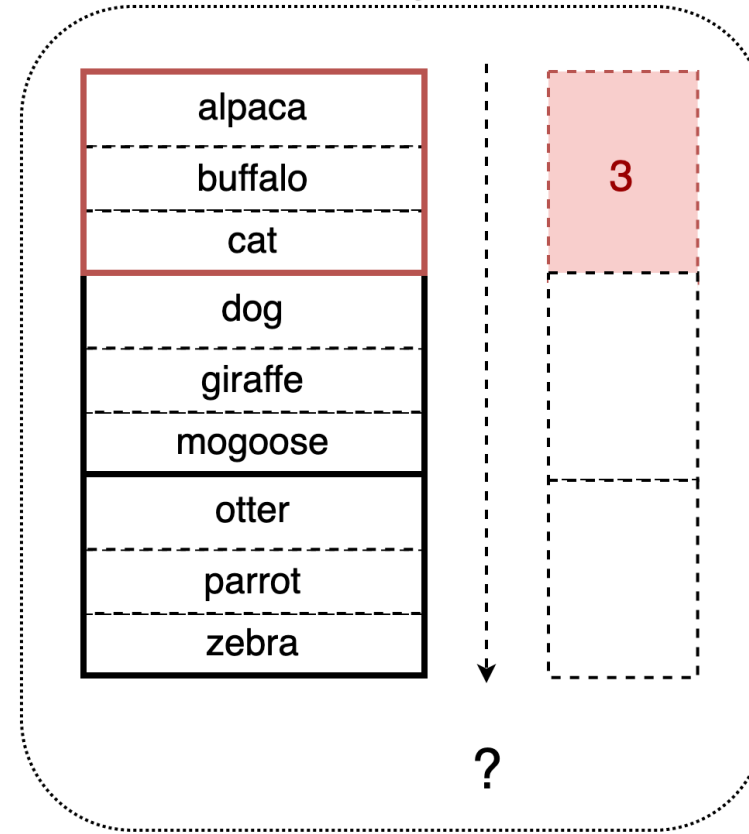


# Побатчевый count

Count by row

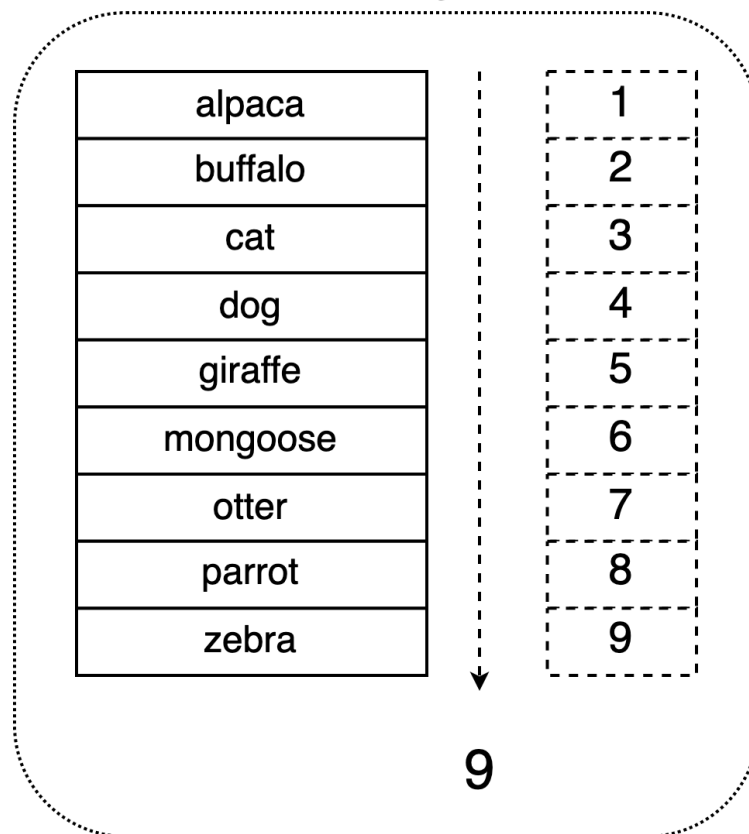


Count by batch

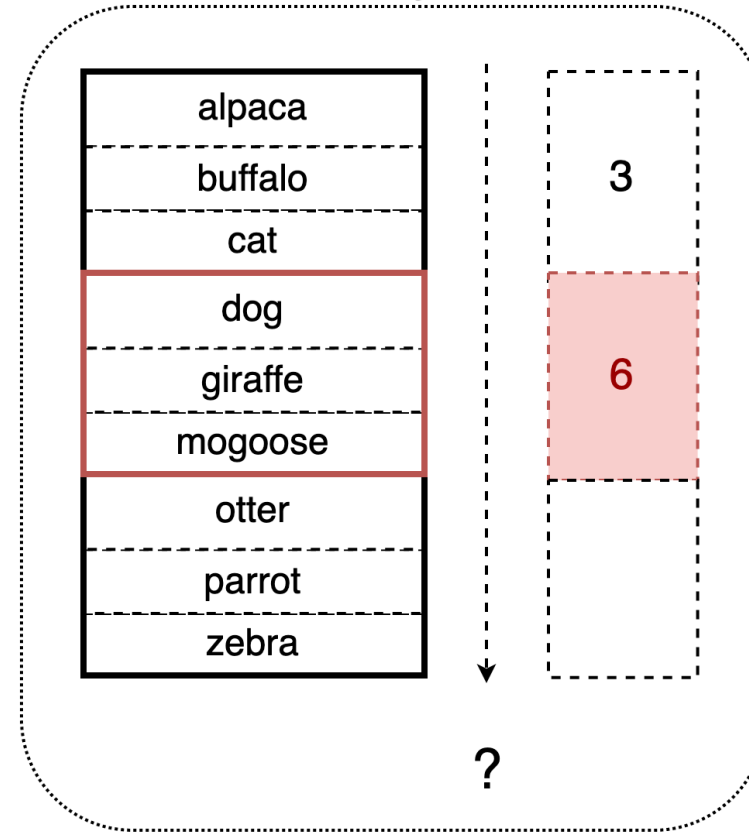


# Побатчевый count

Count by row

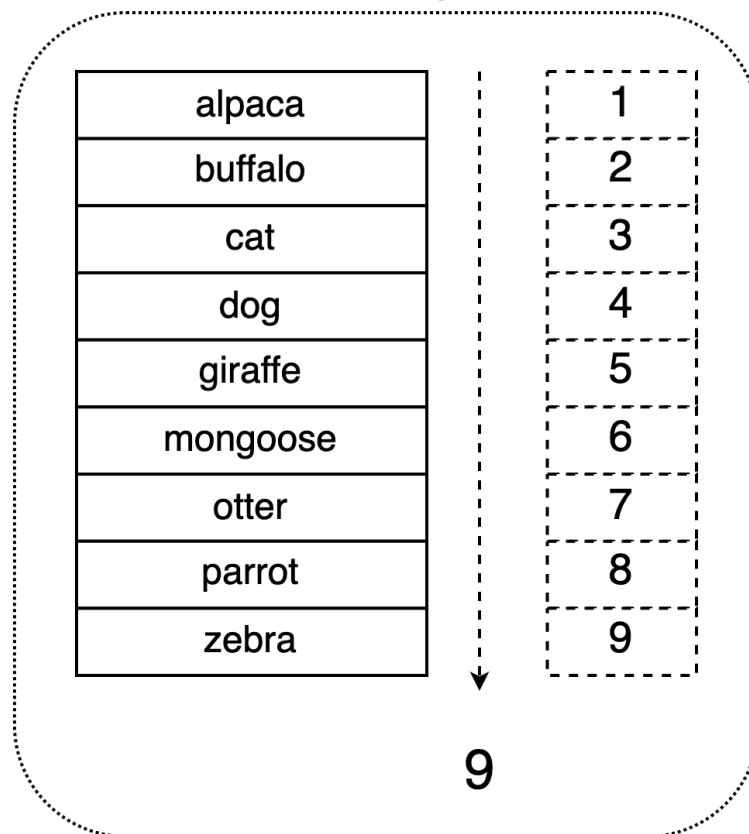


Count by batch

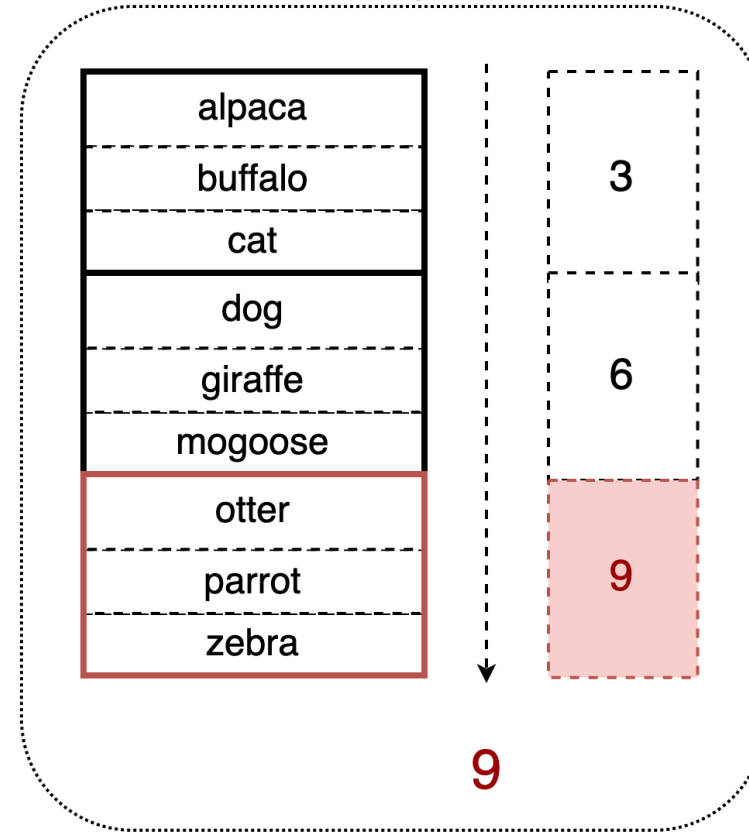


# Побатчевый count

Count by row



Count by batch



# Почему плохо читать построчно

- Не даём спарку использовать метаданные батчей
- Тратим CPU на перекладку из построчного формата в построчный

# Что не так с простым подходом

- Нет чтения батчами
- Нельзя прочитать директорию целиком

# YT – файловая система

## HDFS

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> dataset_1
|           |-----> part-0.parquet
|           |-----> part-1.parquet
|           |-----> part-2.parquet
|       |---> dataset_2
|           |-----> part-0.parquet
|           |-----> part-1.parquet
```

## YT

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |-----> example_table_1
|       |-----> example_table_2
```

# Чтение директории из HDFS

## HDFS

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> dataset_1
|           |-----> part-0.parquet
|           |-----> part-1.parquet
|           |-----> part-2.parquet
|       |---> dataset_2
|           |-----> part-0.parquet
|           |-----> part-1.parquet
```

## YT

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |-----> example_table_1
|       |-----> example_table_2
```

```
spark.read.parquet("/tmp/example")
```

# Чтение директории из YT

## HDFS

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> dataset_1
|           |-----> part-0.parquet
|           |-----> part-1.parquet
|           |-----> part-2.parquet
|       |---> dataset_2
|           |-----> part-0.parquet
|           |-----> part-1.parquet
```

```
spark.read.parquet("/tmp/example")
```

## YT

```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |-----> example_table_1
|       |-----> example_table_2
```

```
spark.read.yt("/tmp/example")
```



# Partition discovery

HDFS

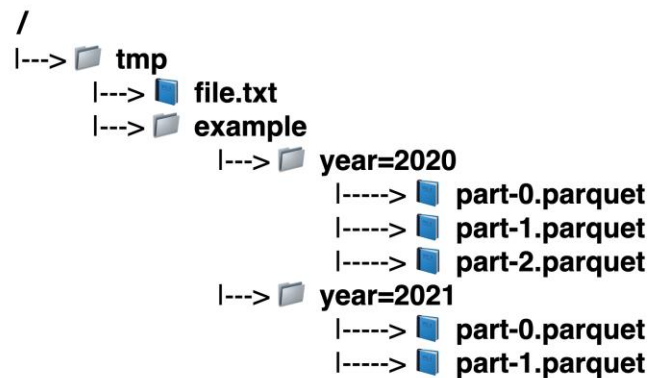
```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> year=2020
|           |---> part-0.parquet
|           |---> part-1.parquet
|           |---> part-2.parquet
|       |---> year=2021
|           |---> part-0.parquet
|           |---> part-1.parquet
```

```
spark.read.parquet("/tmp/example")
```

a	..	year
cat	..	2020
dog	..	2020
otter	..	2021

# Partition discovery

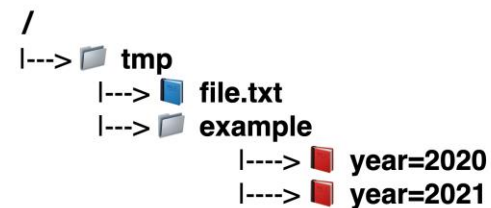
HDFS



```
spark.read.parquet("/tmp/example")
```

a	..	year
cat	..	2020
dog	..	2020
otter	..	2021

YT



```
spark.read.yt("/tmp/example")
```

a	..	year
cat	..	2020
dog	..	2020
otter	..	2021

# Чего мы хотим от нового подхода

- Научиться читать поколоночными батчами
- Переиспользовать код, написанный для HDFS и Parquet:
  - Partition discovery
  - Обход директорий на чтении
  - Создание / удаление временных файлов на записи

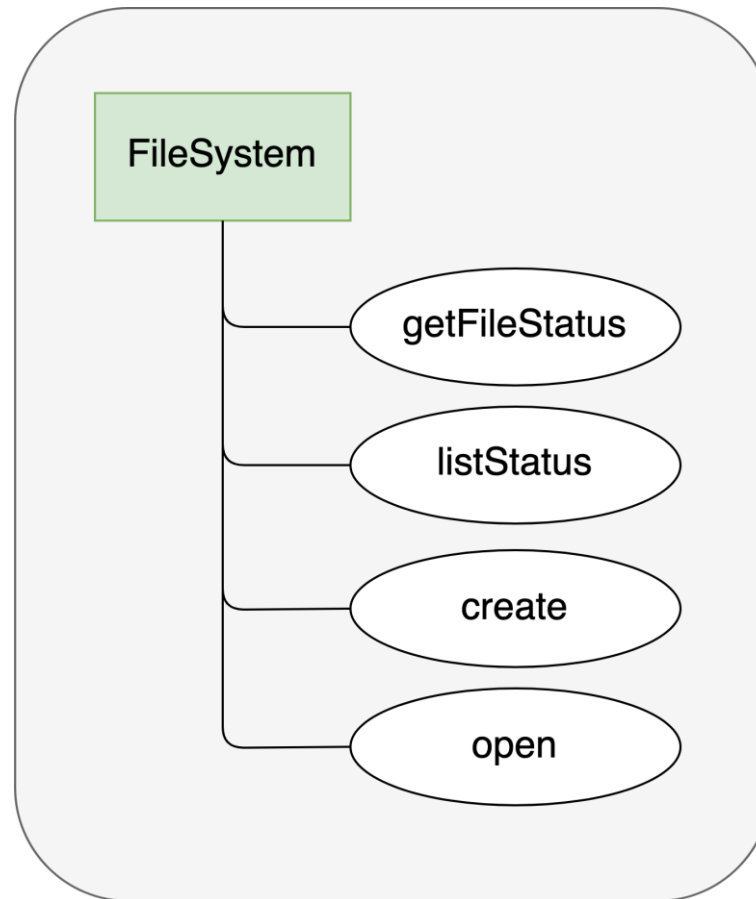
# Новое решение

# Как Spark читает Parquet с HDFS

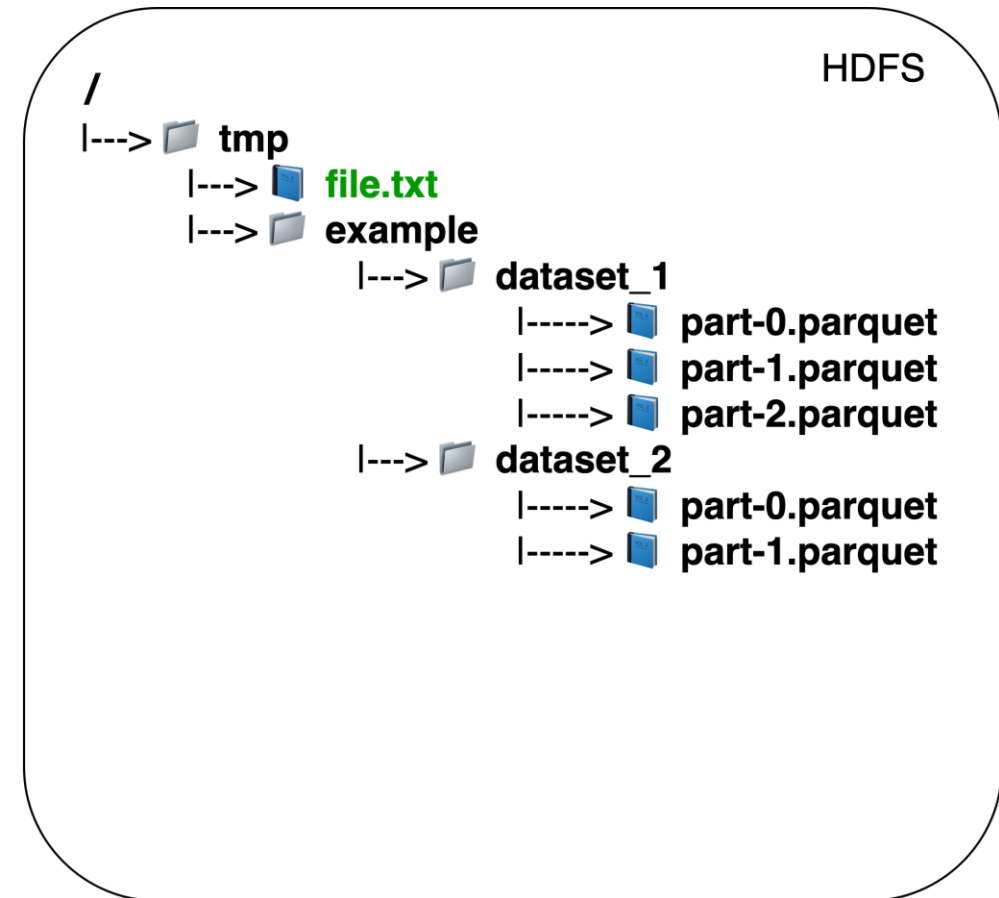
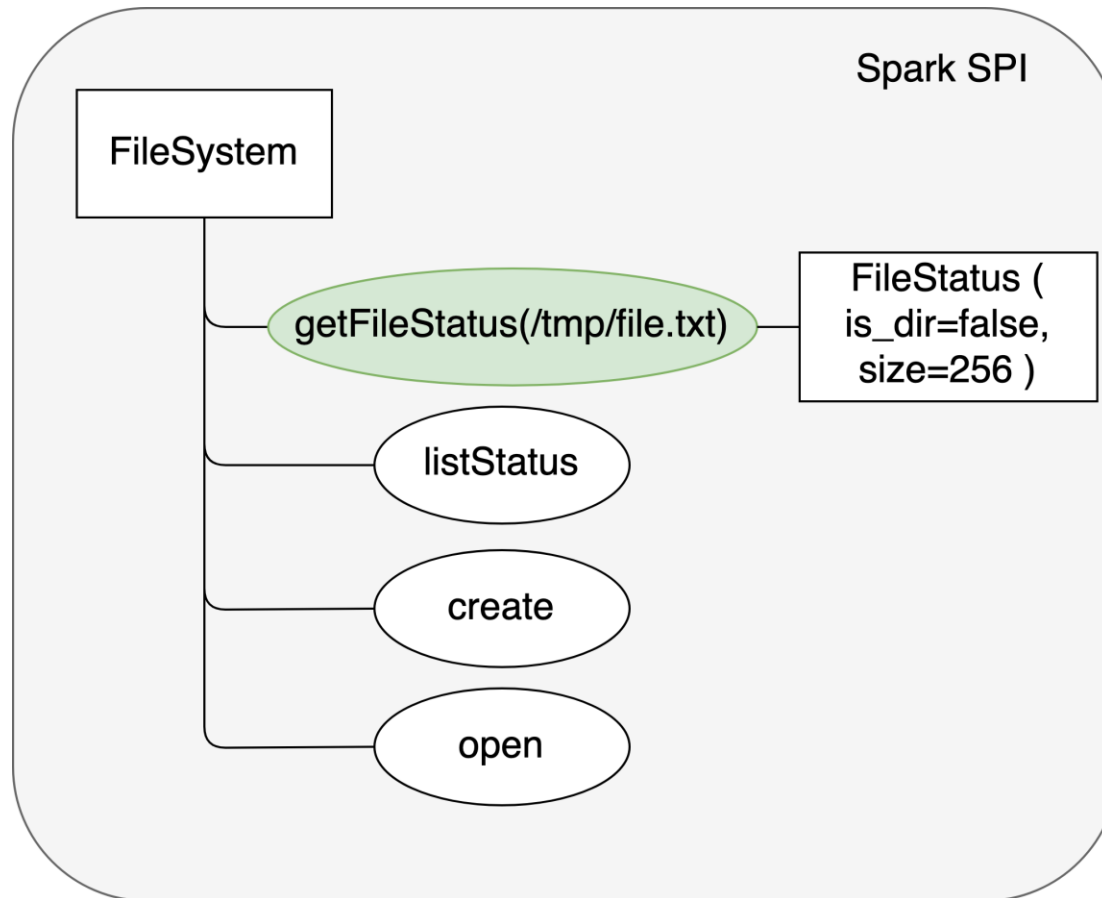
# Разделение DataSource на две части

- Данные лежат в файловой системе (например, HDFS)
- Данные сохранены в каком-то формате (например, Parquet)

# FileSystem

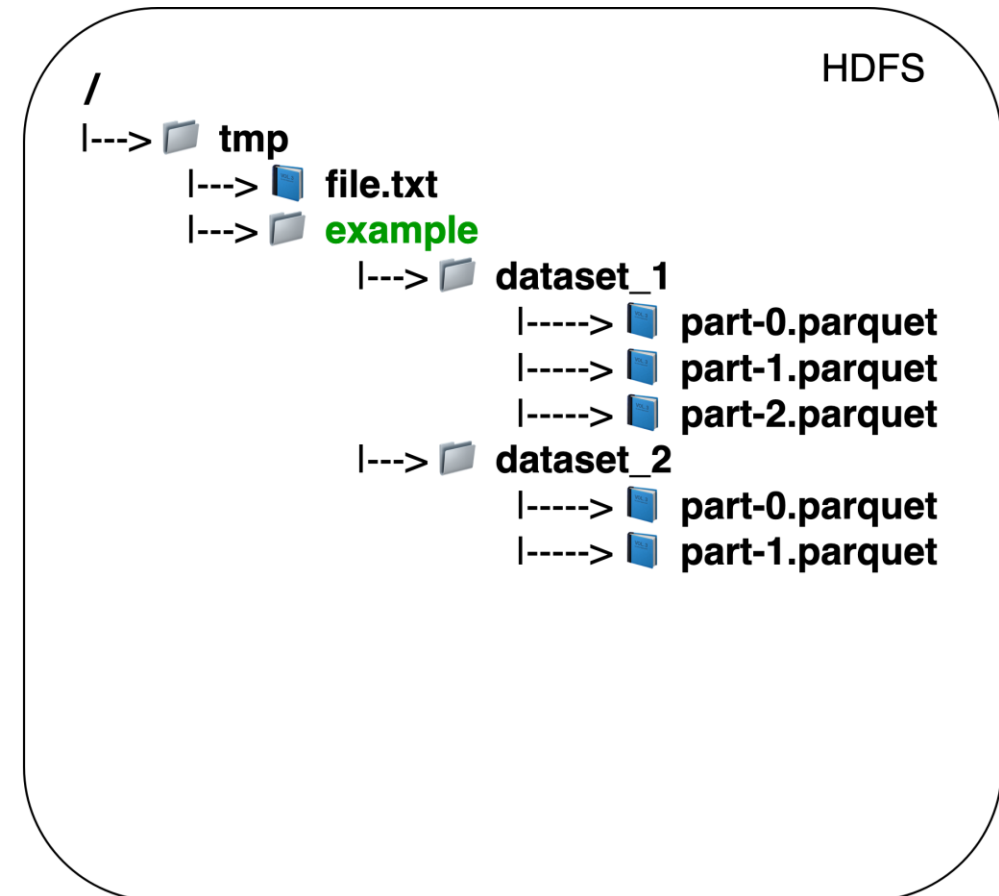
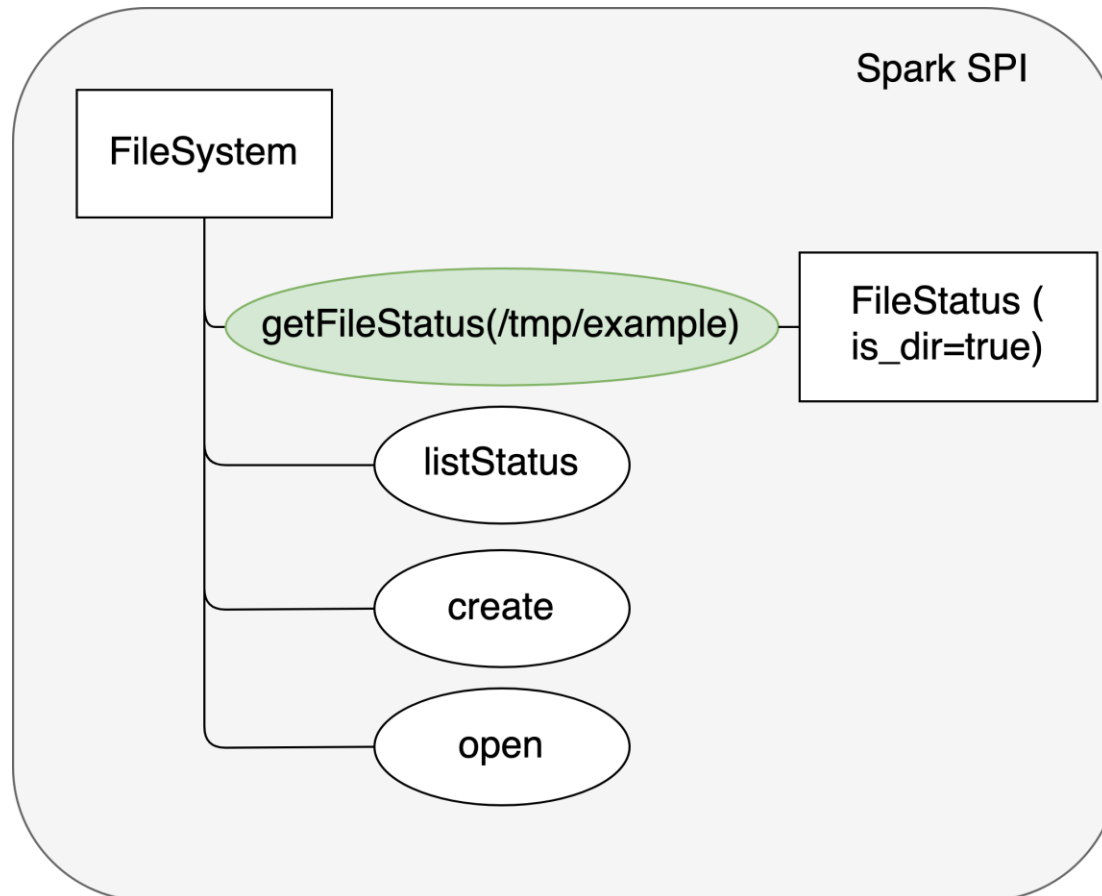


# Статус файла

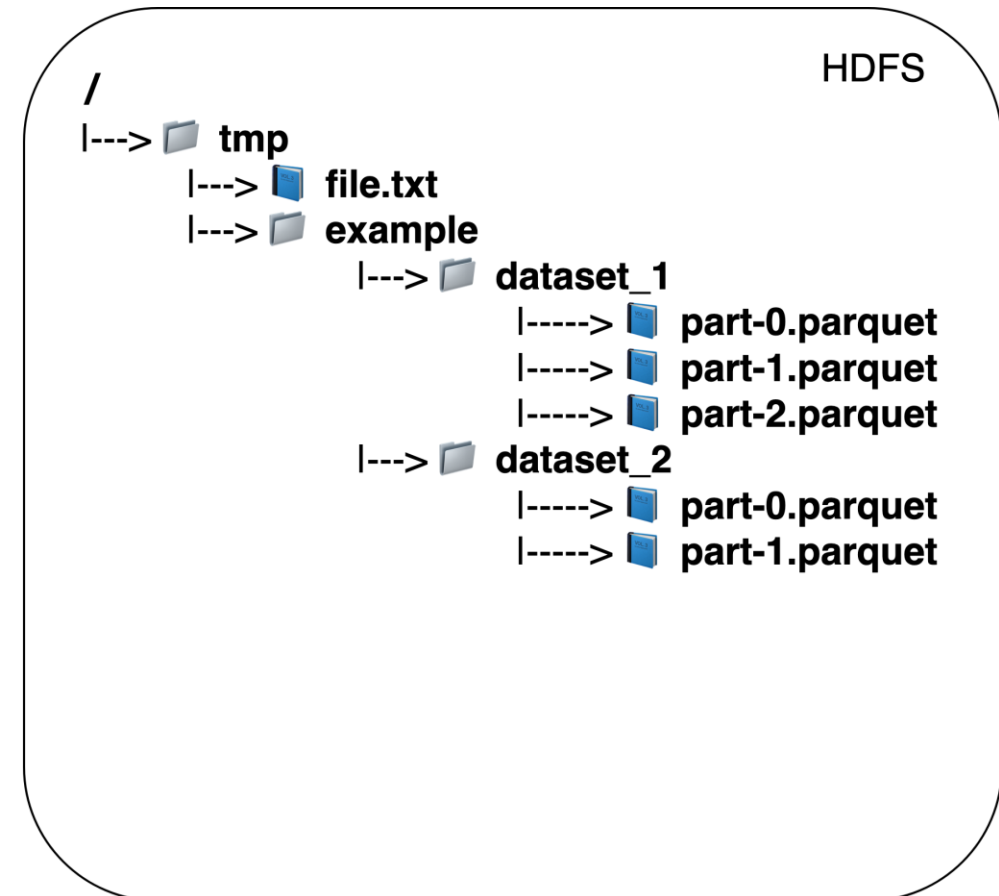
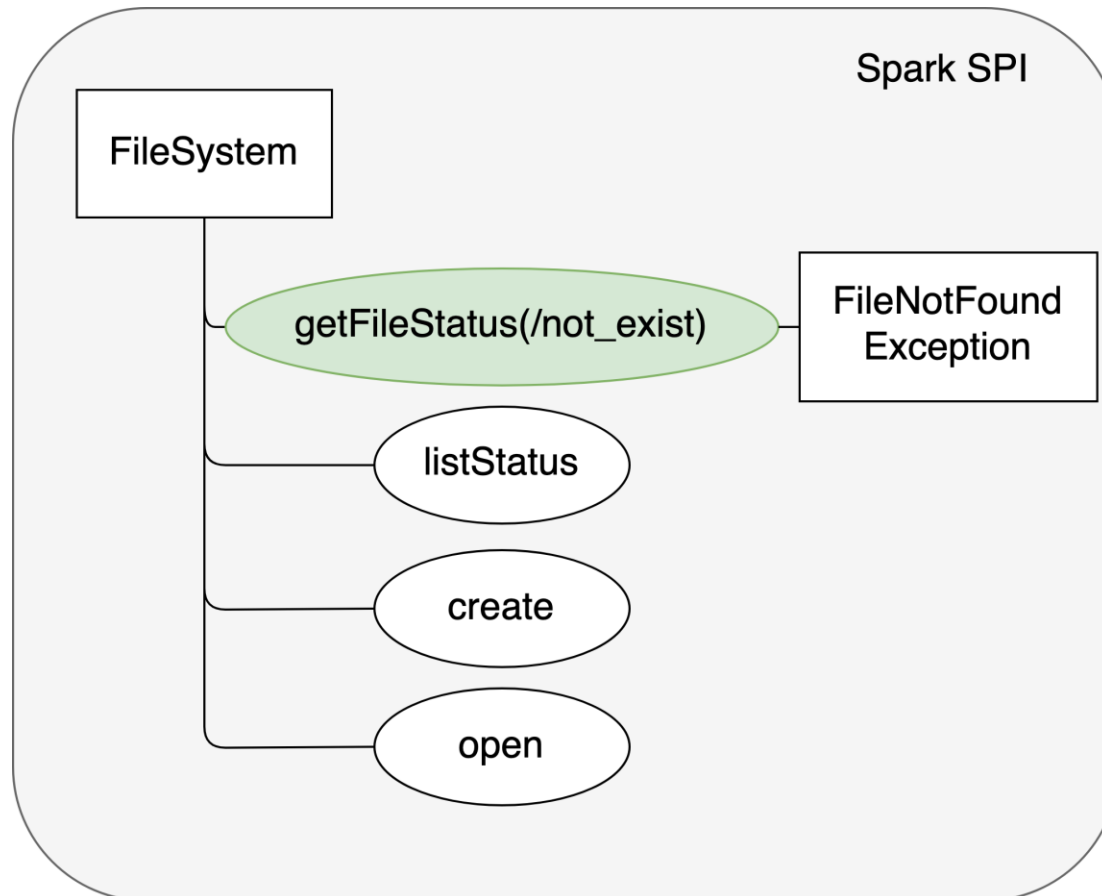




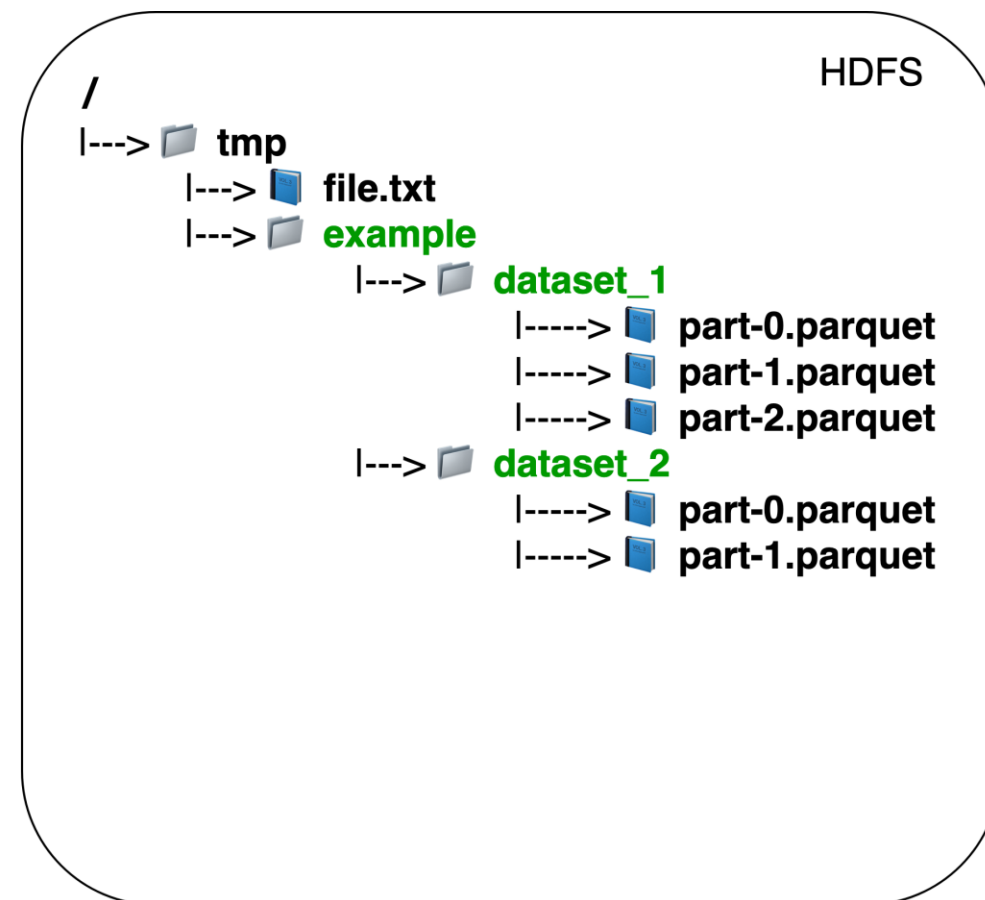
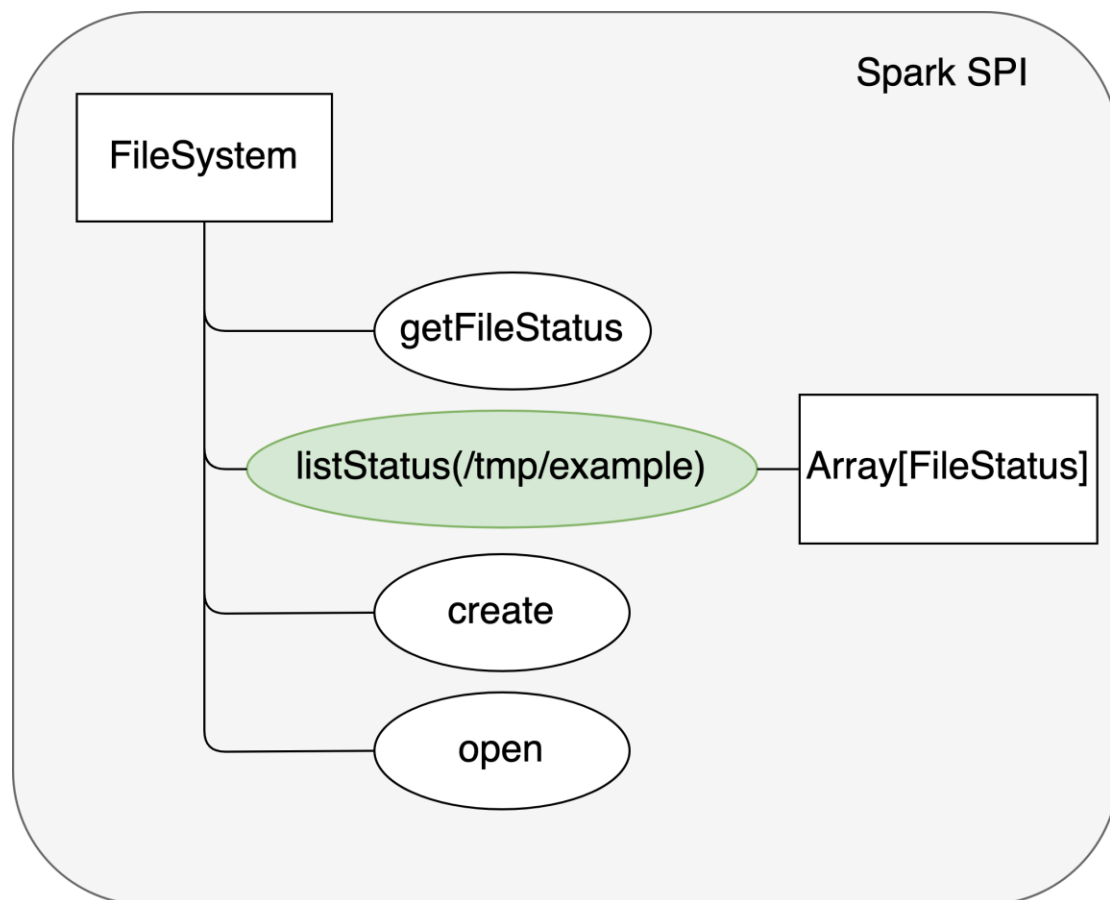
# Статус директории



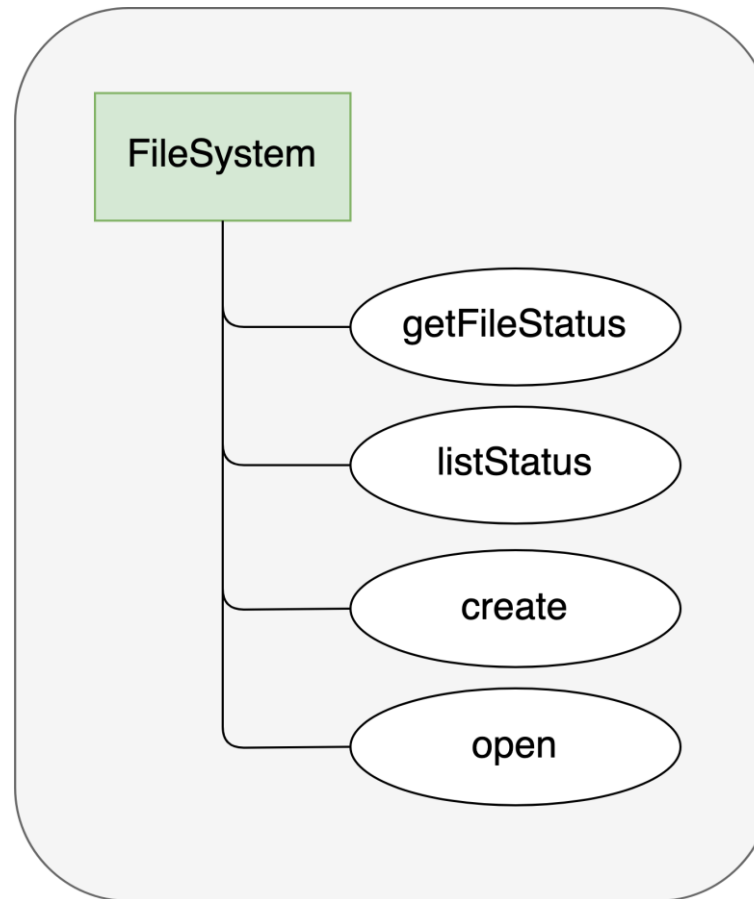
# Файл не существует



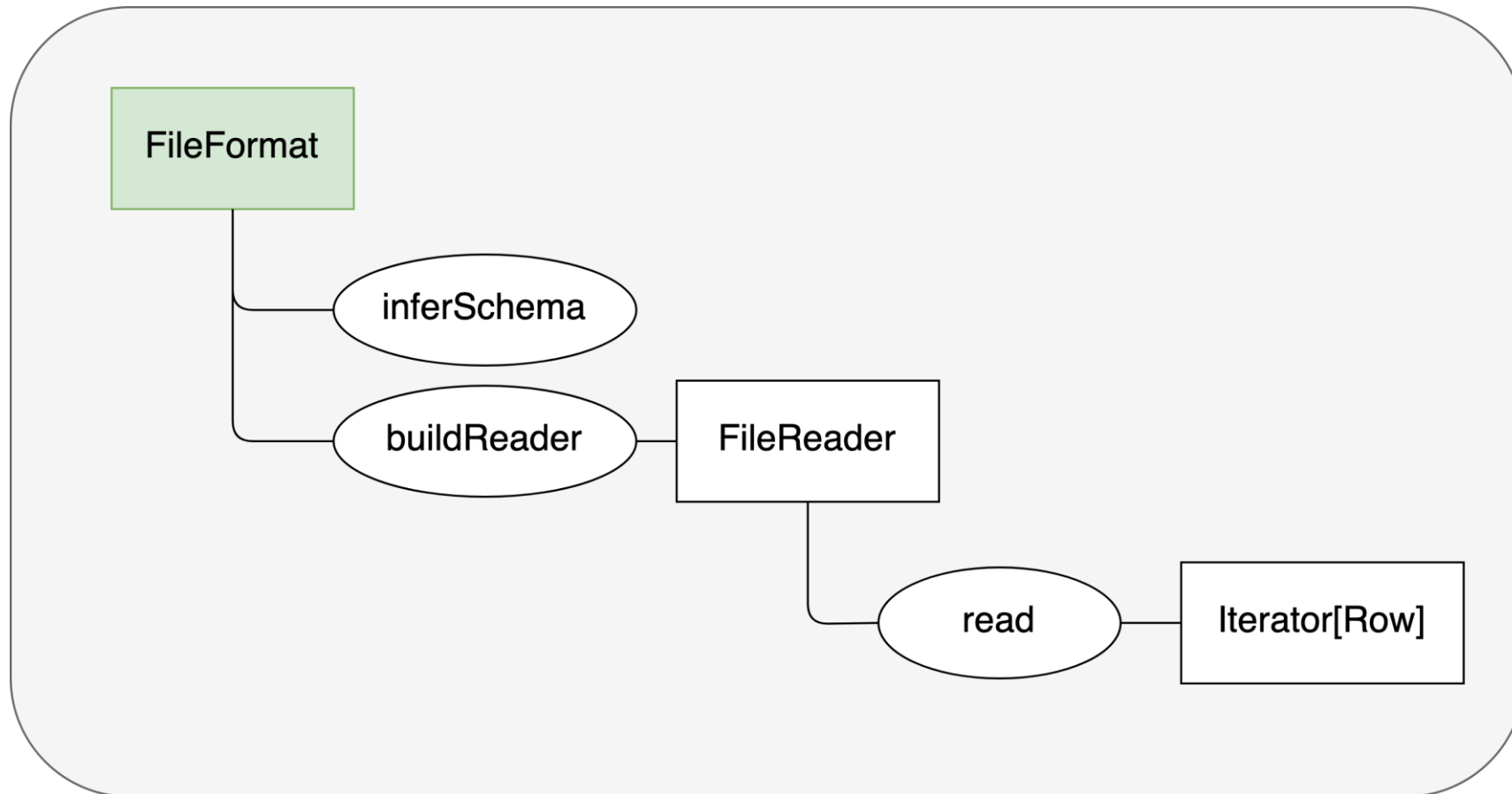
# Листинг директории



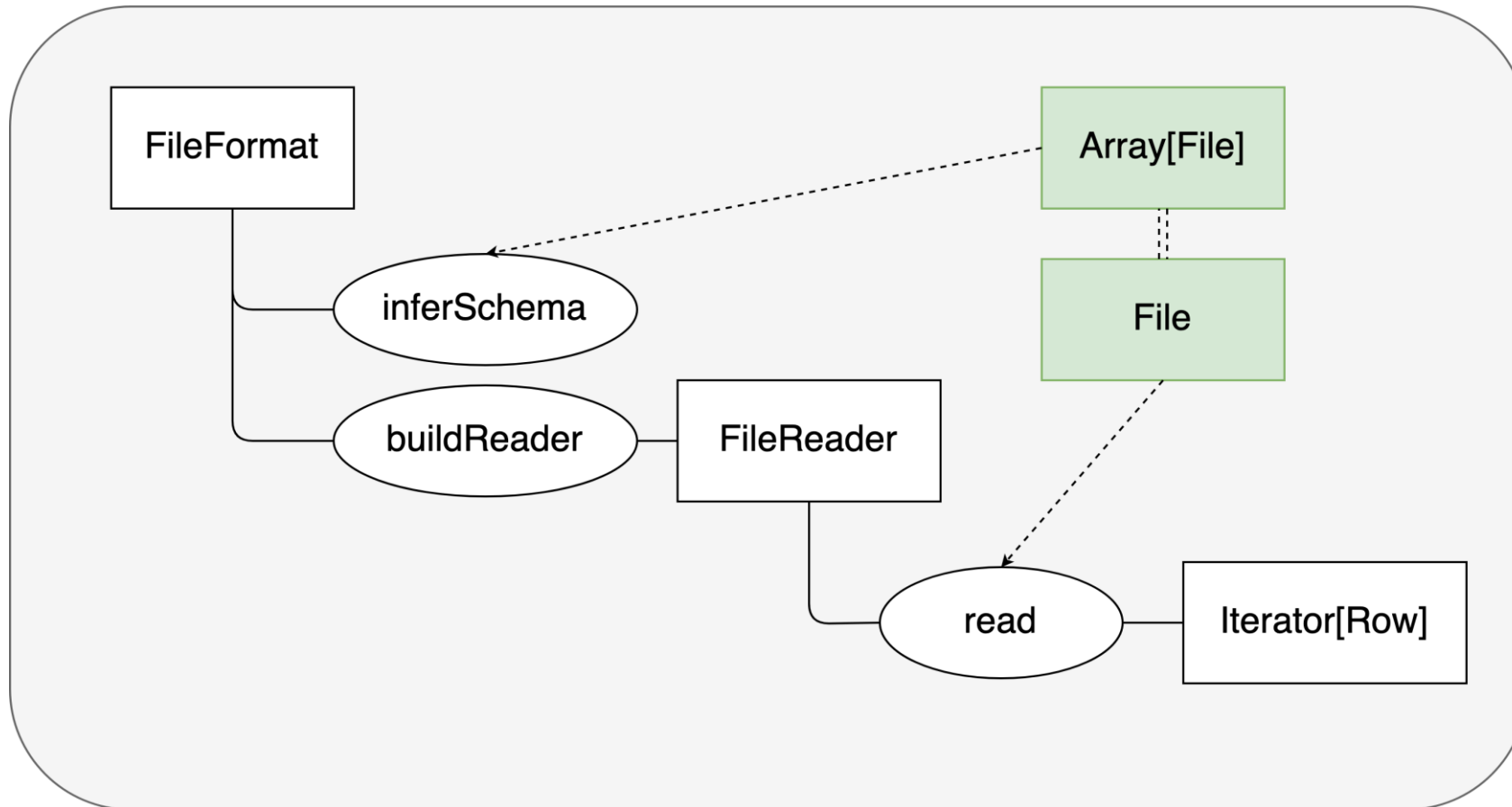
# FileSystem



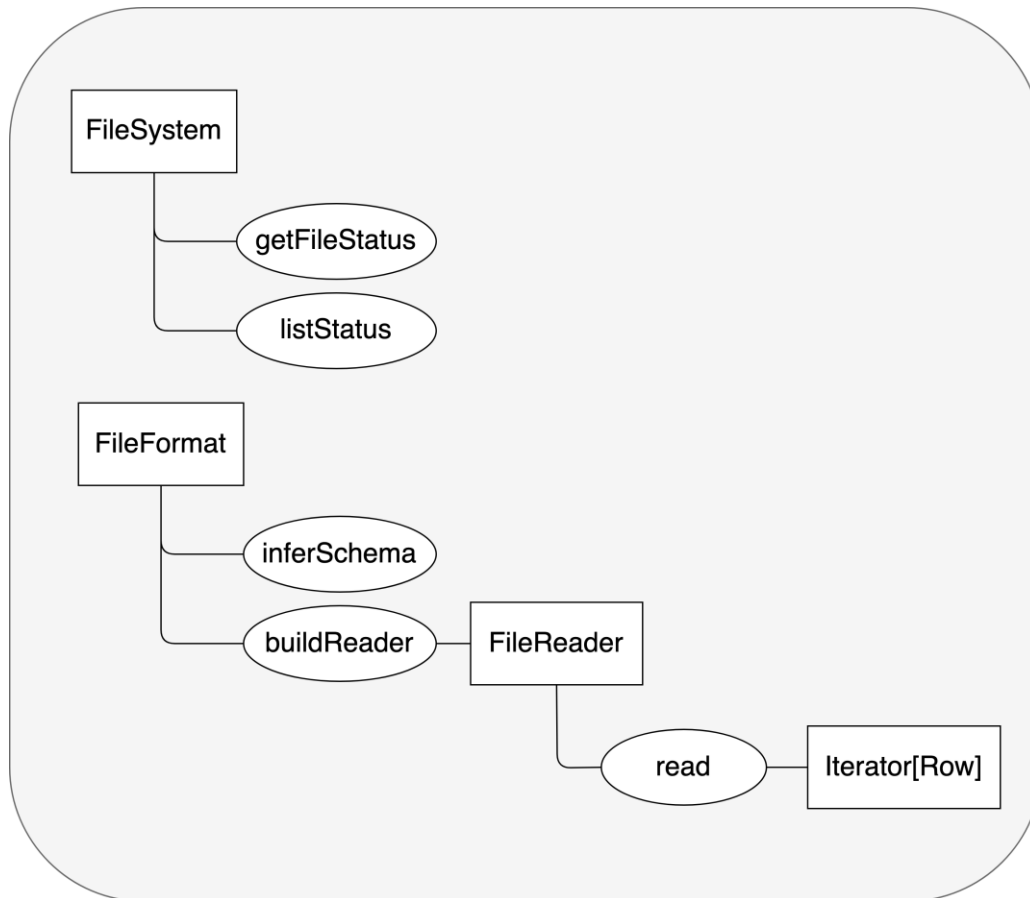
# FileFormat



# FileFormat

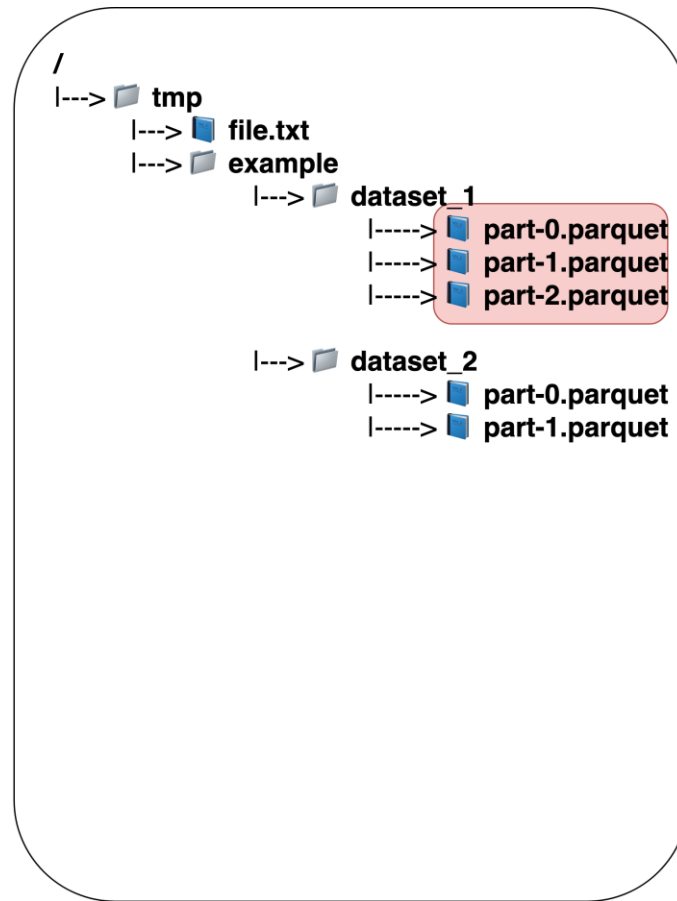
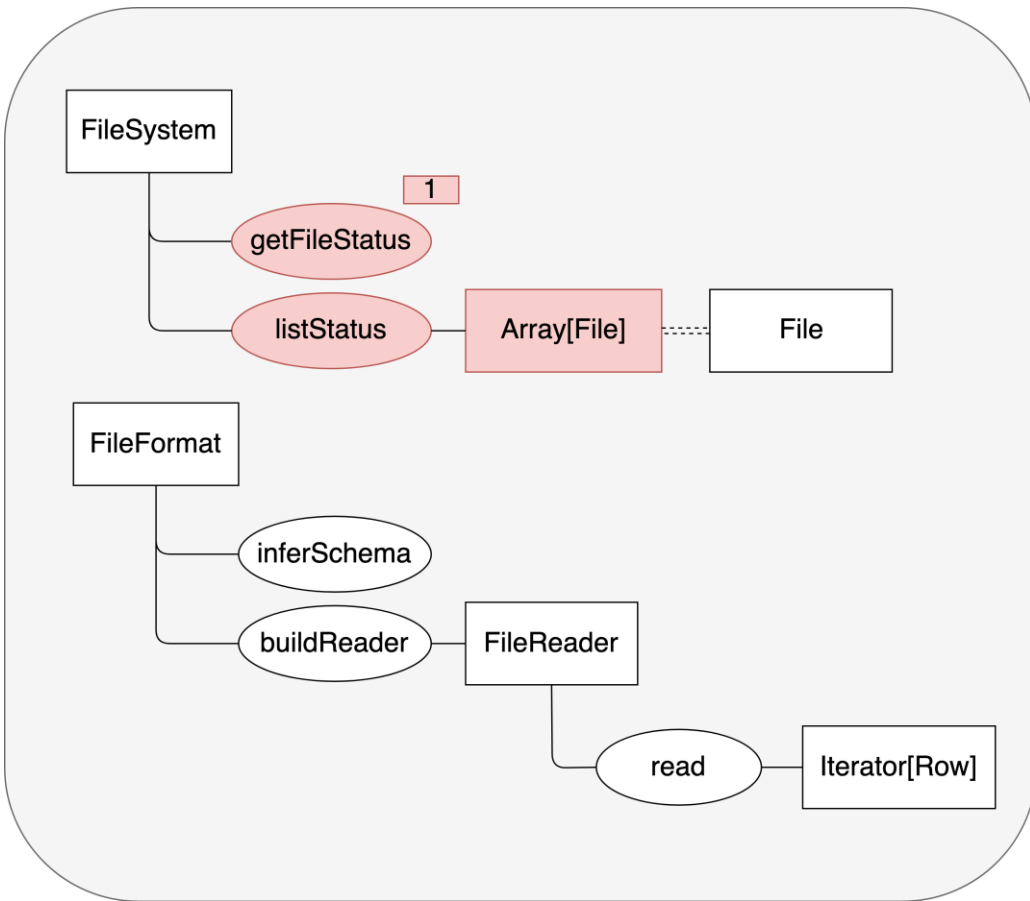


# Чтение с HDFS



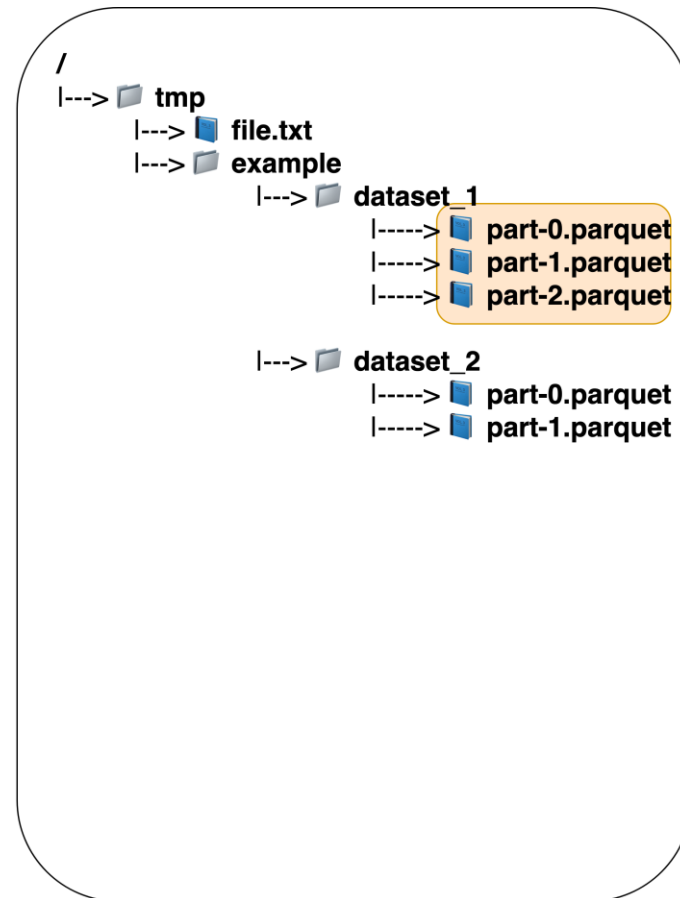
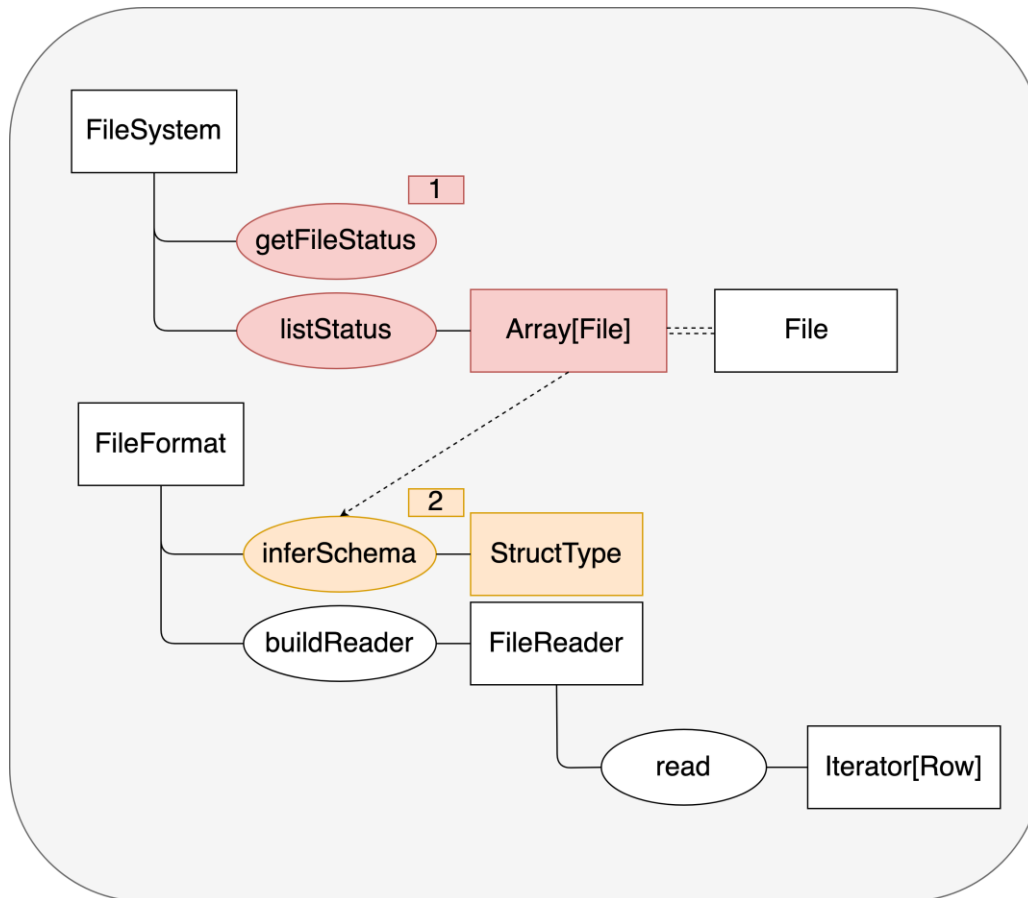
```
/
|---> tmp
|   |---> file.txt
|   |---> example
|       |---> dataset_1
|           |---> part-0.parquet
|           |---> part-1.parquet
|           |---> part-2.parquet
|       |---> dataset_2
|           |---> part-0.parquet
|           |---> part-1.parquet
```

# Листинг директории

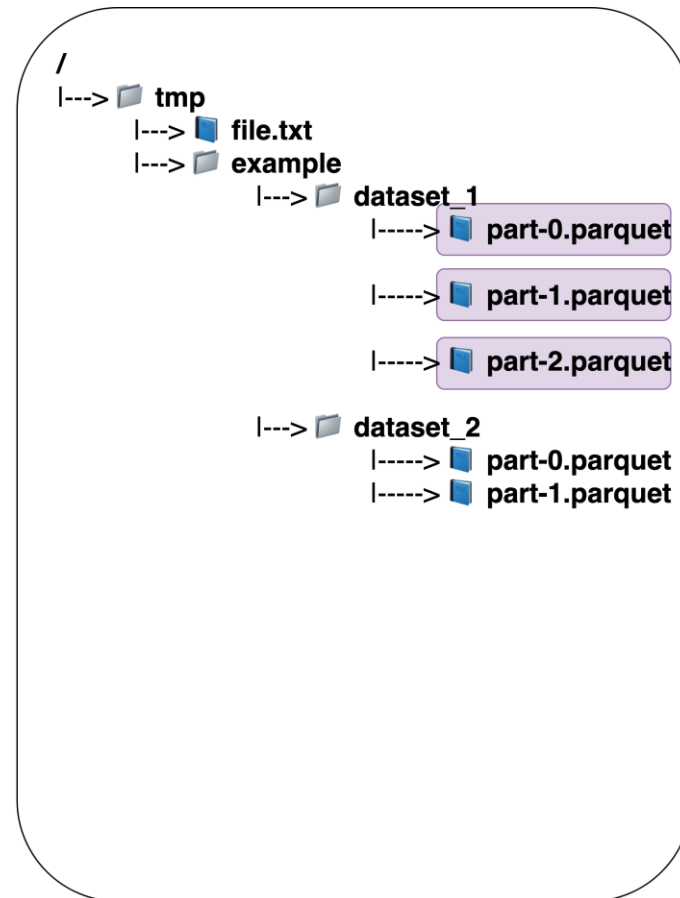
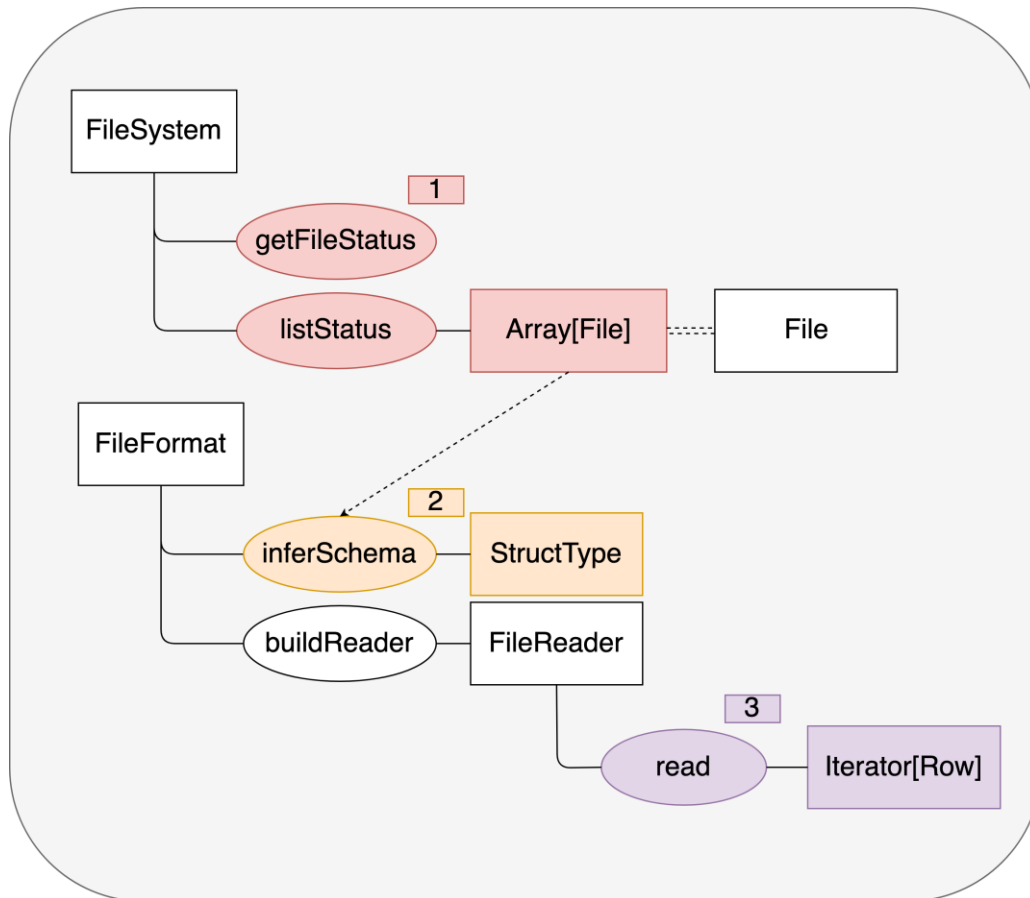




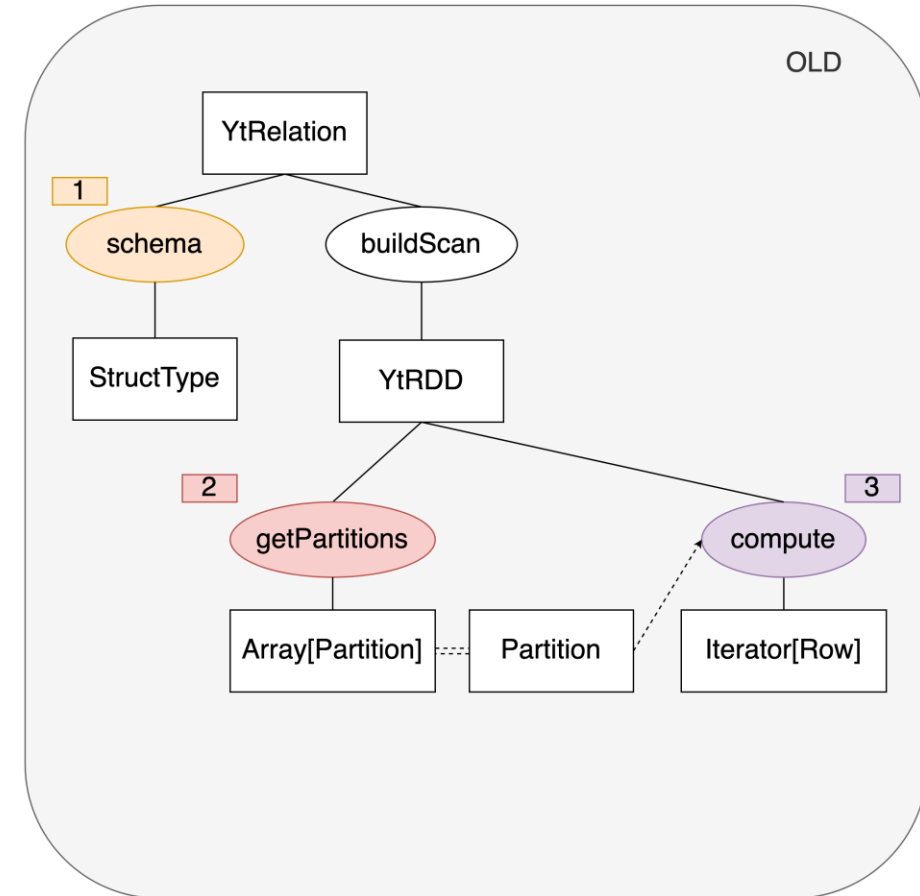
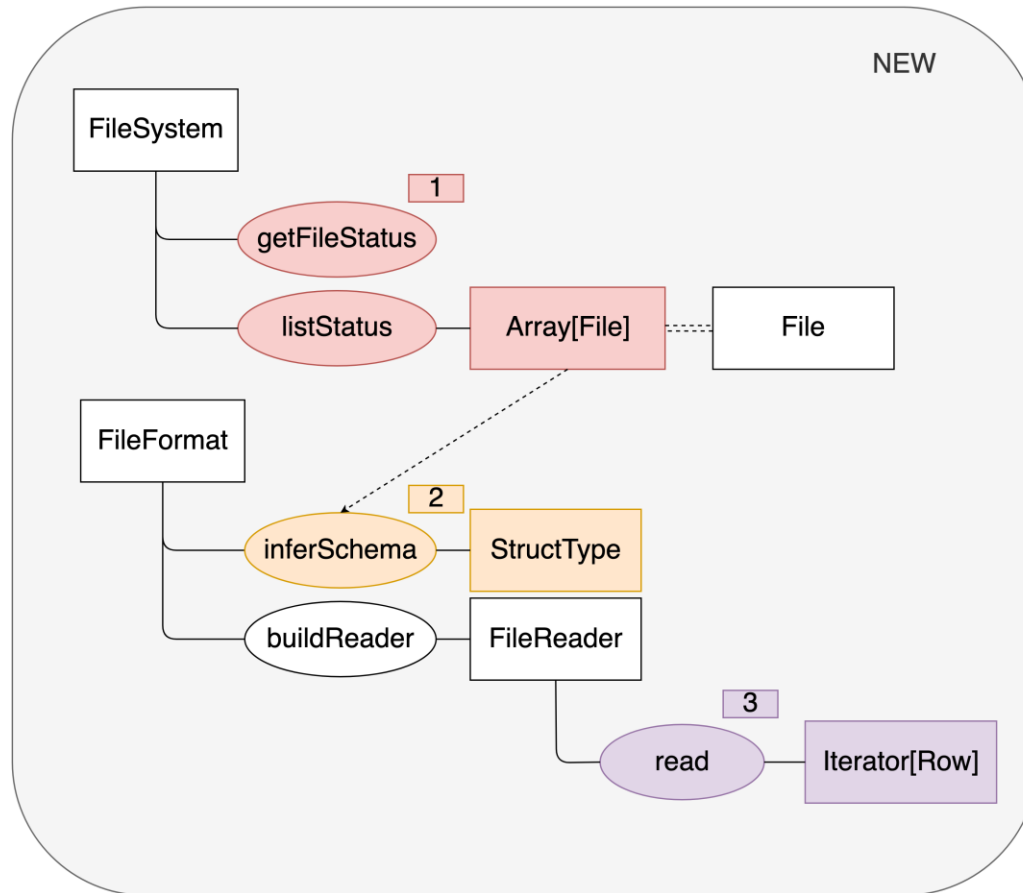
# Чтение схемы



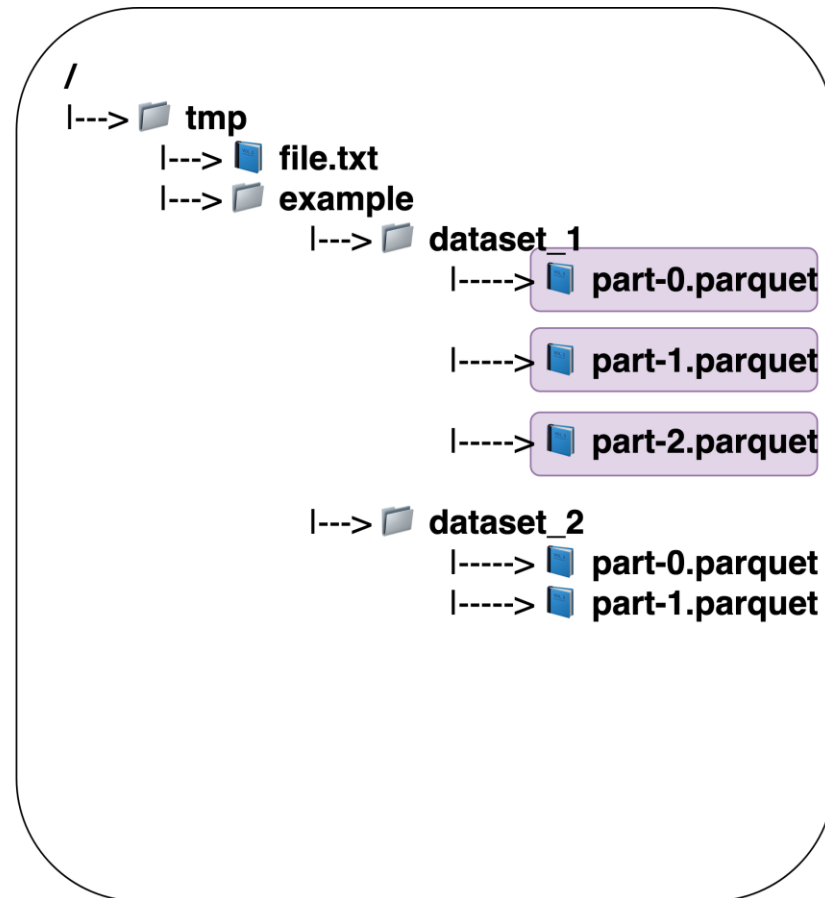
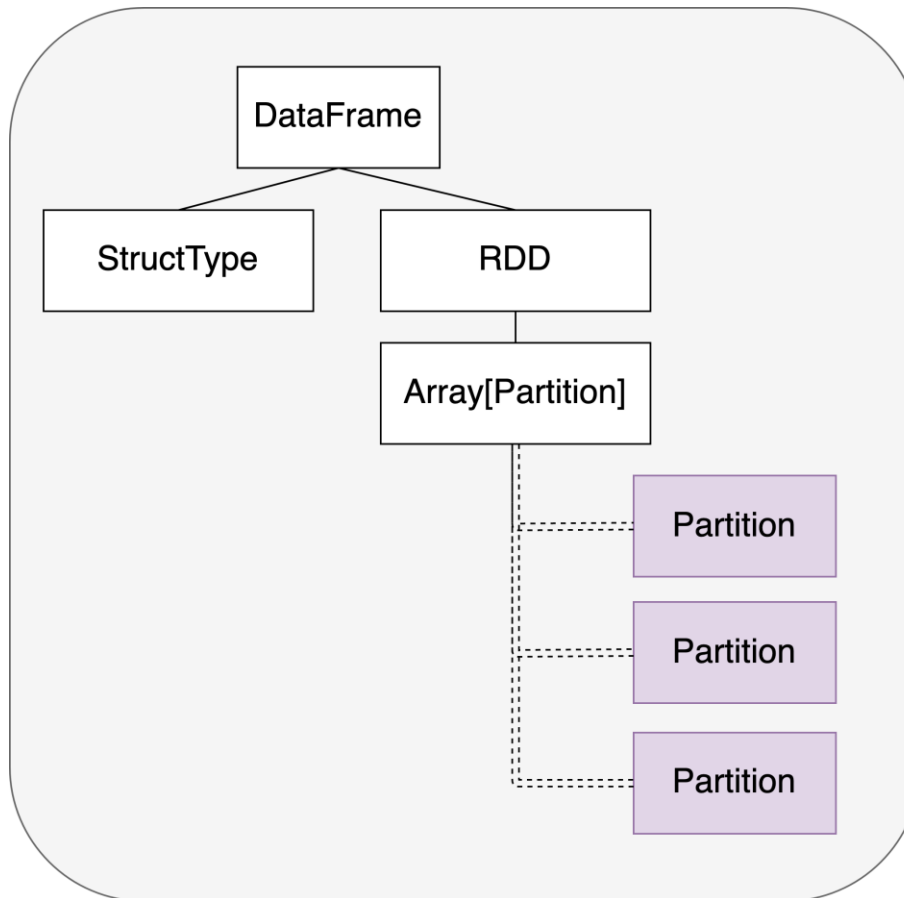
# Чтение партии



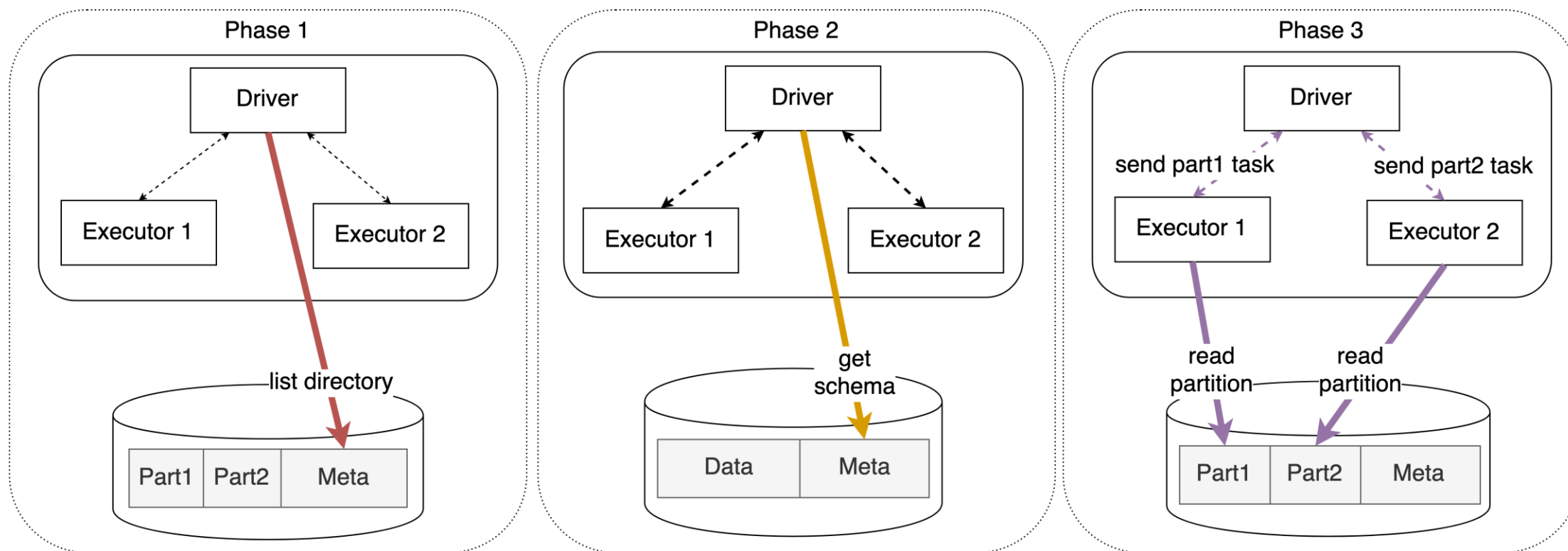
# Сравнение с простым подходом



# Файл = партиция

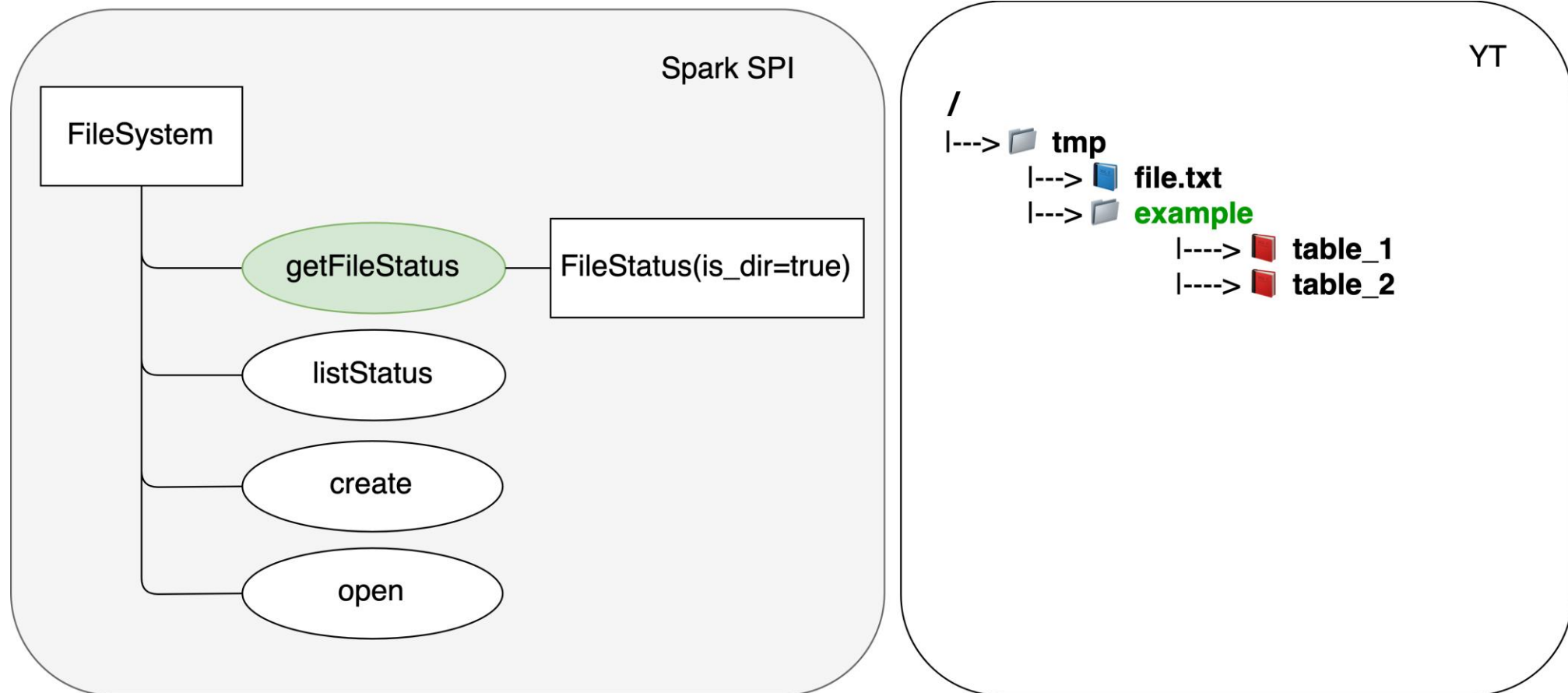


# Как это выполняется

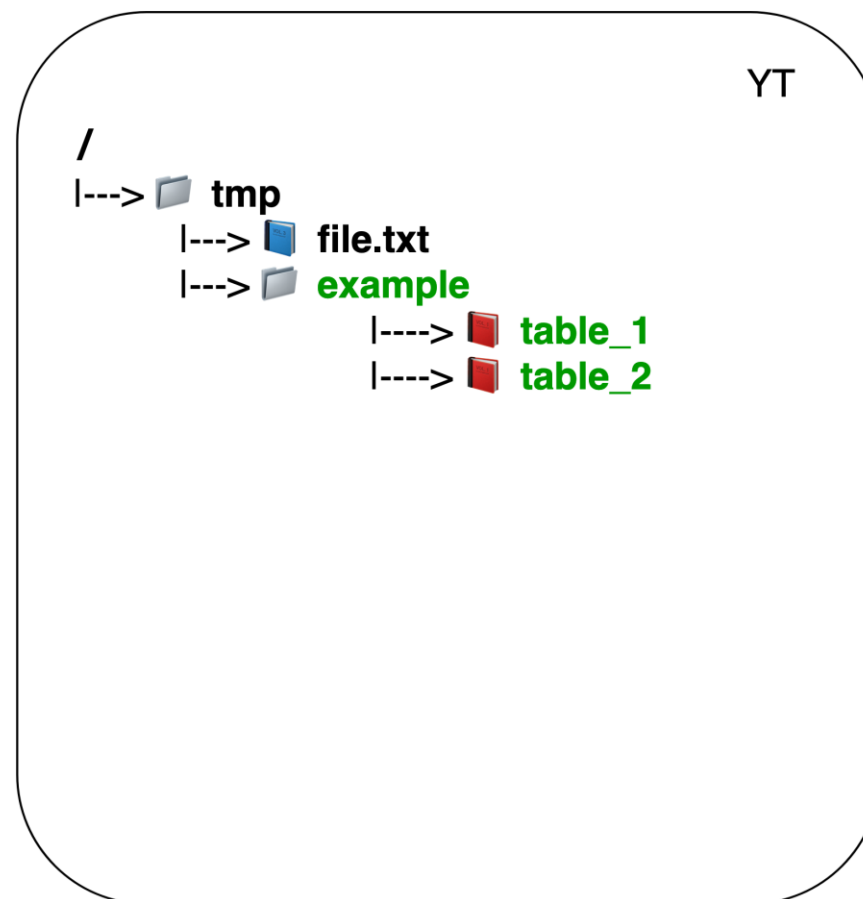
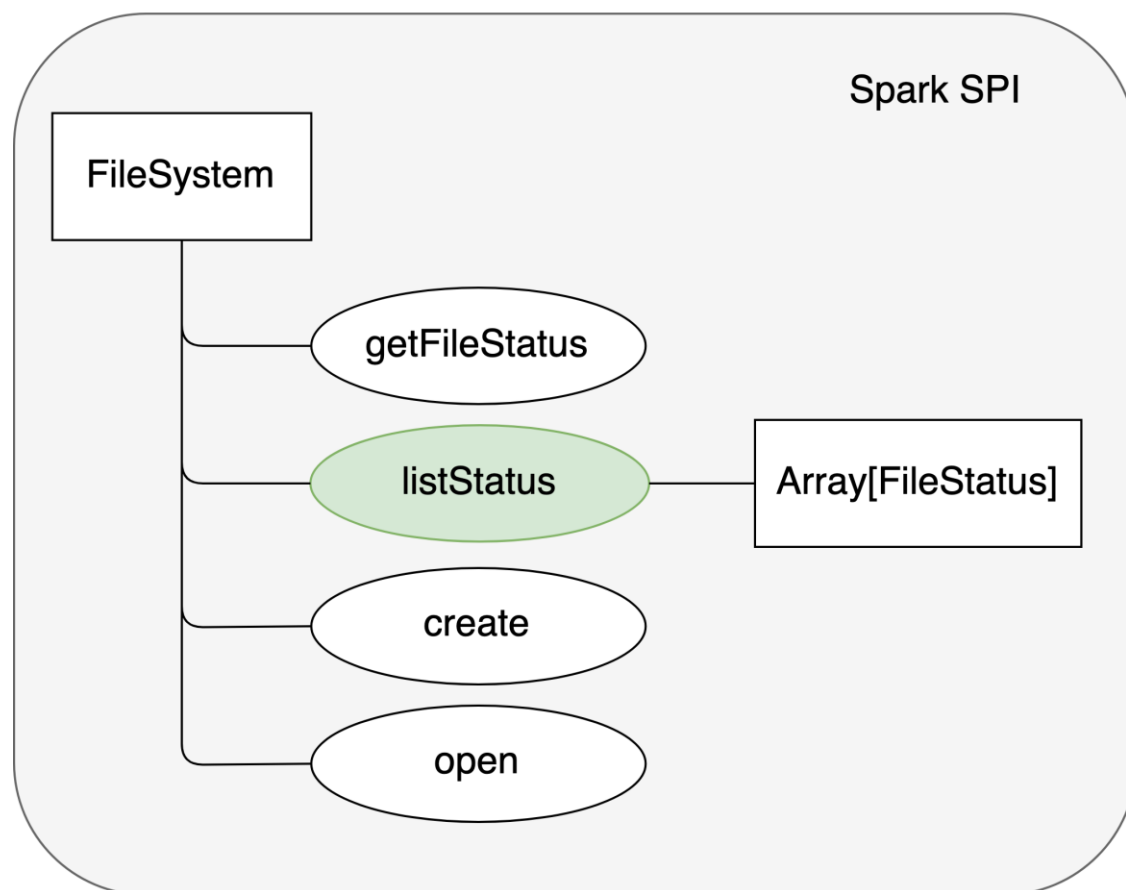


# Как SPYT читает таблицы с YT

# YtFileSystem

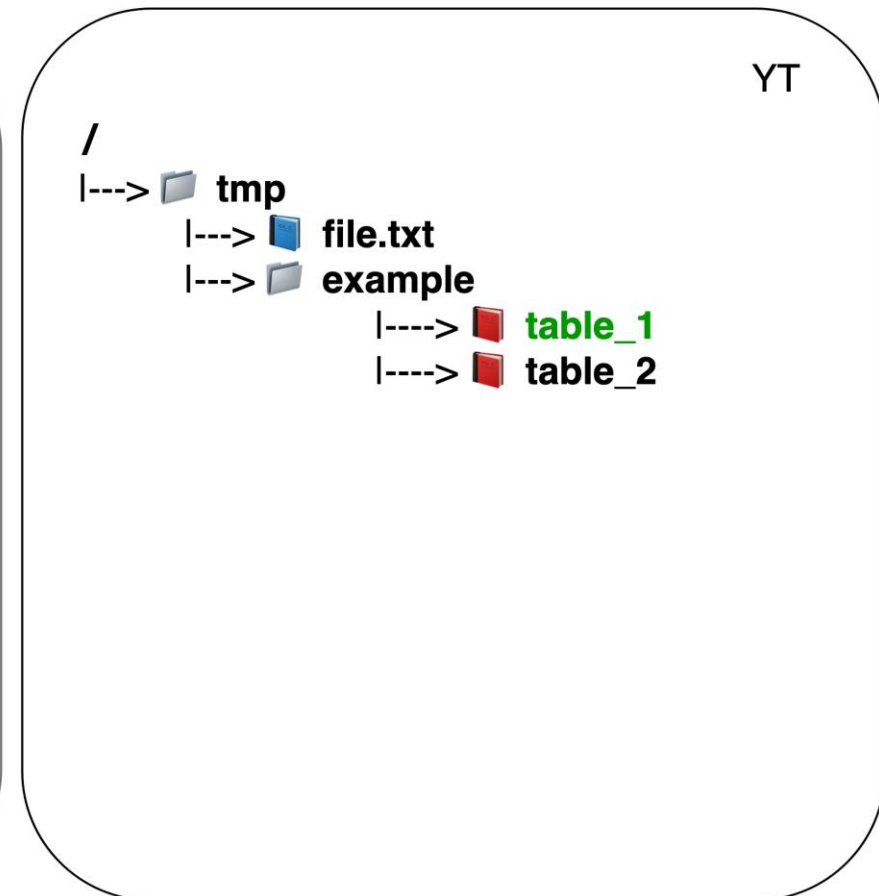
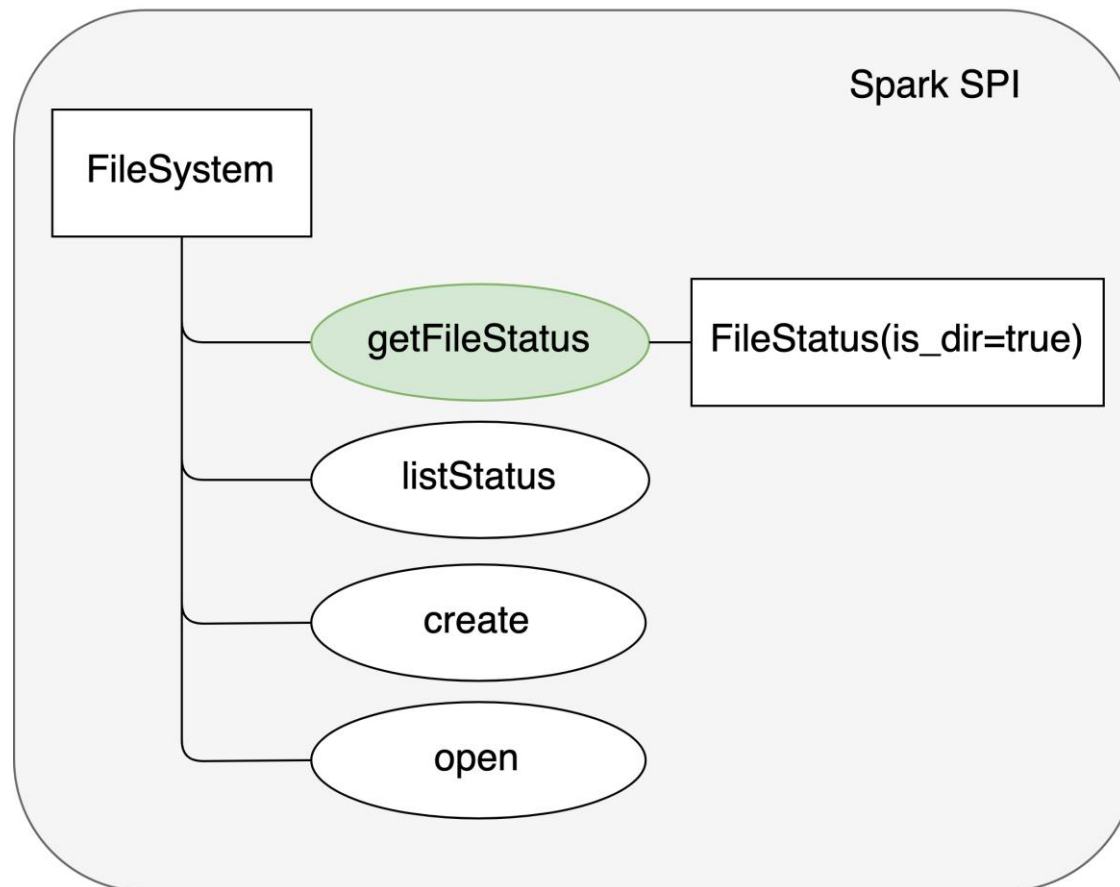


# Листинг директории

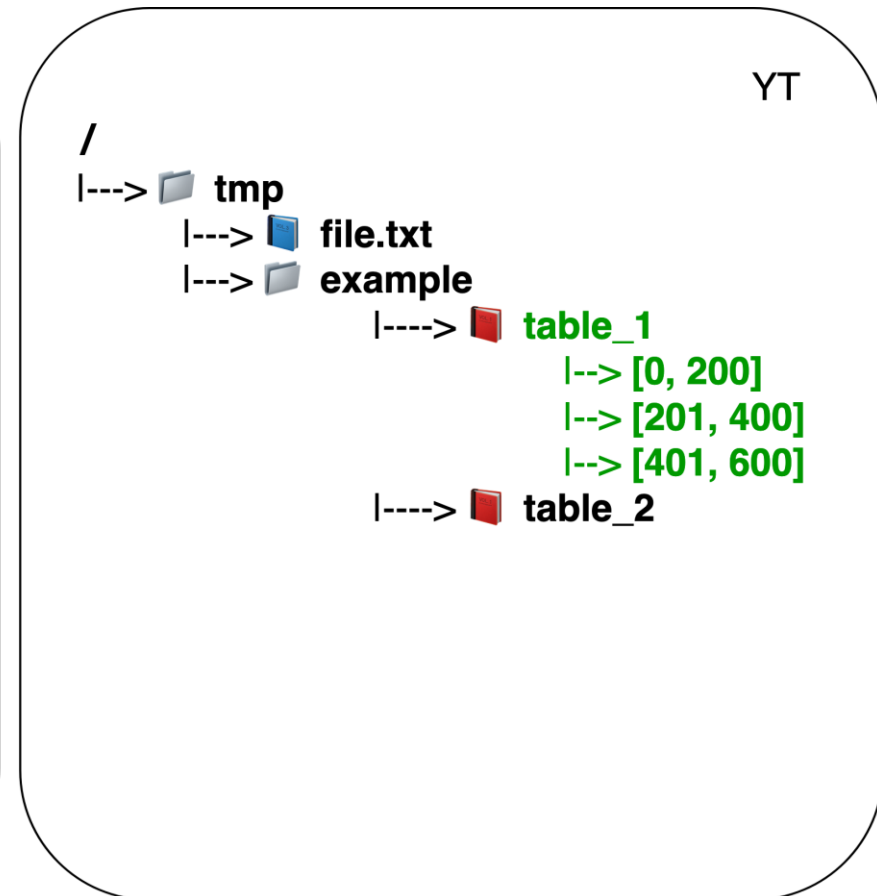
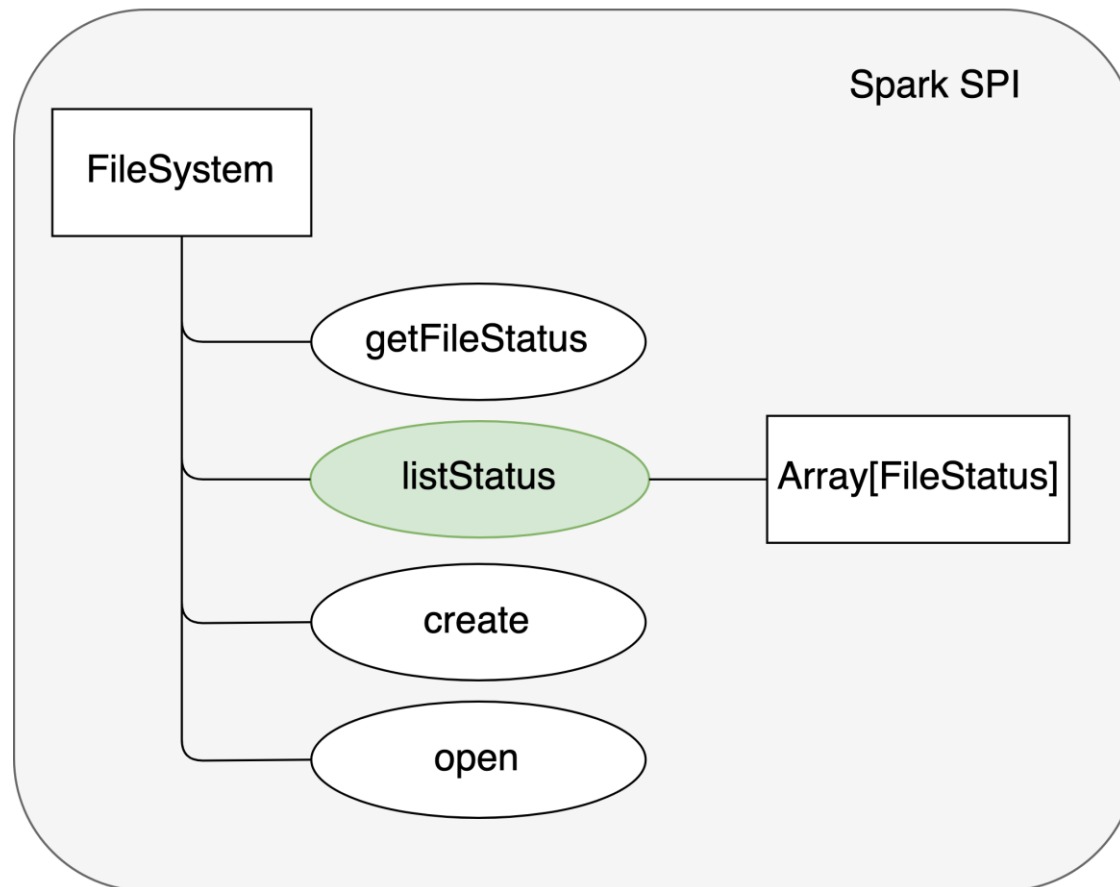




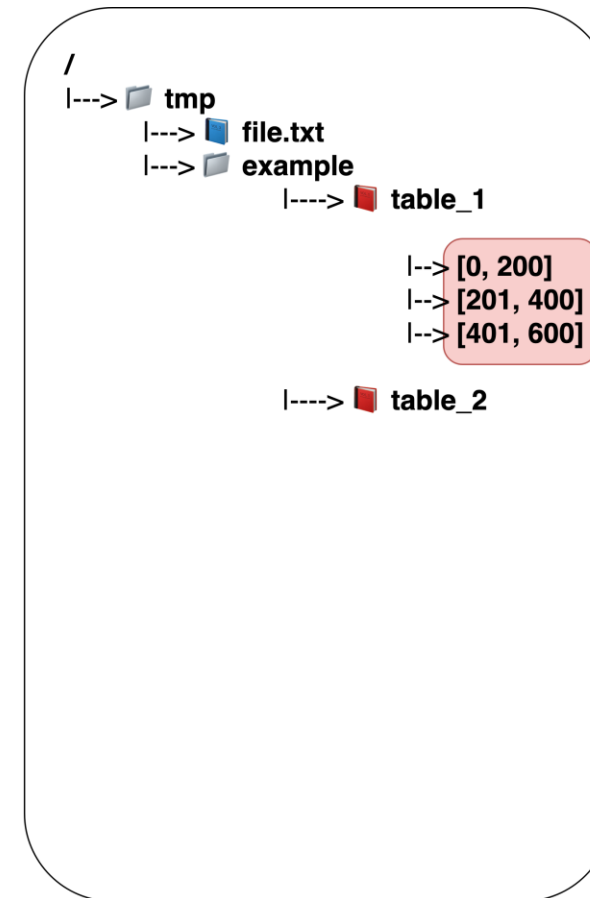
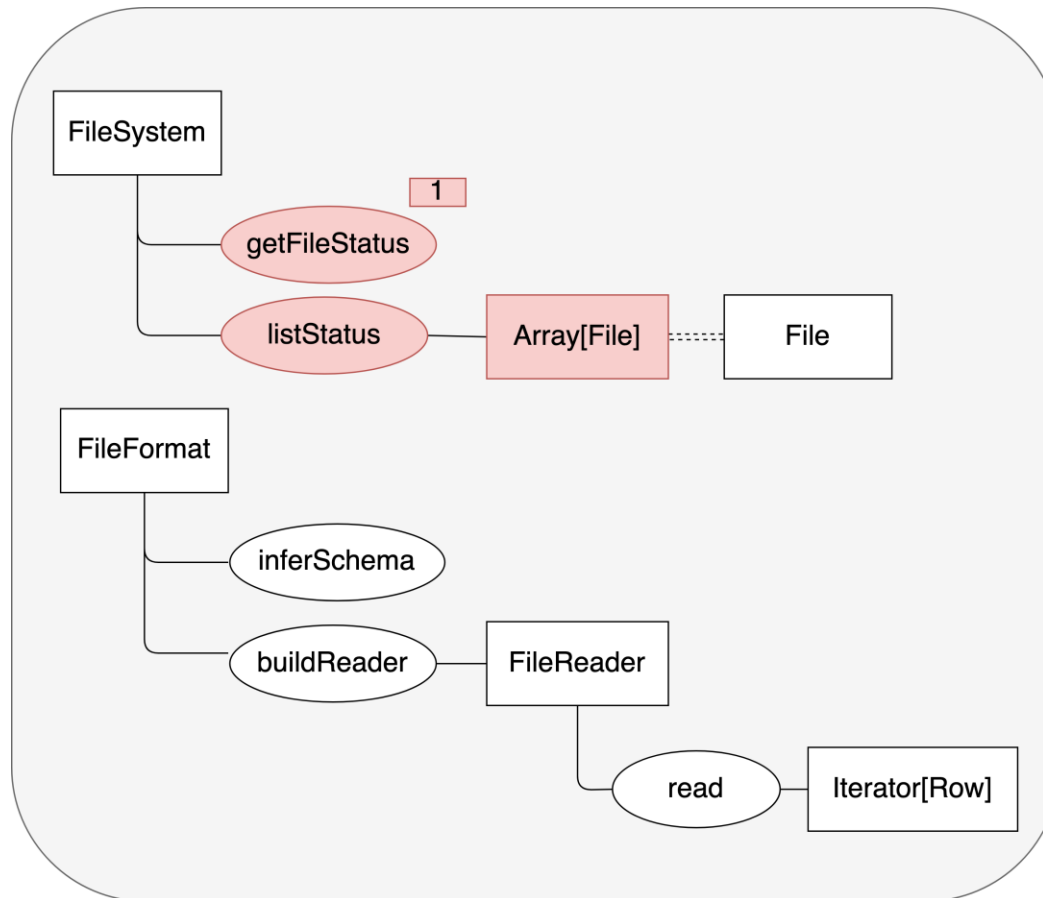
# Таблица = директория



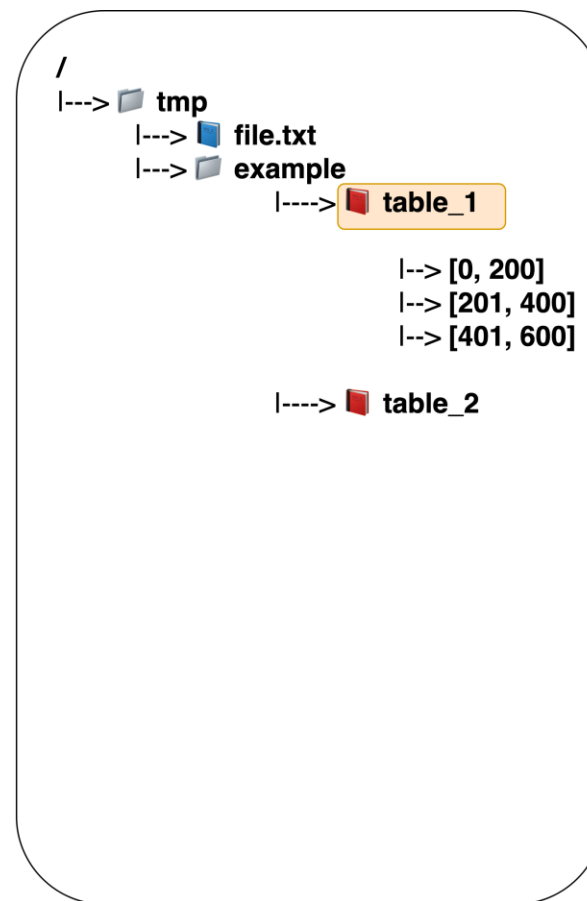
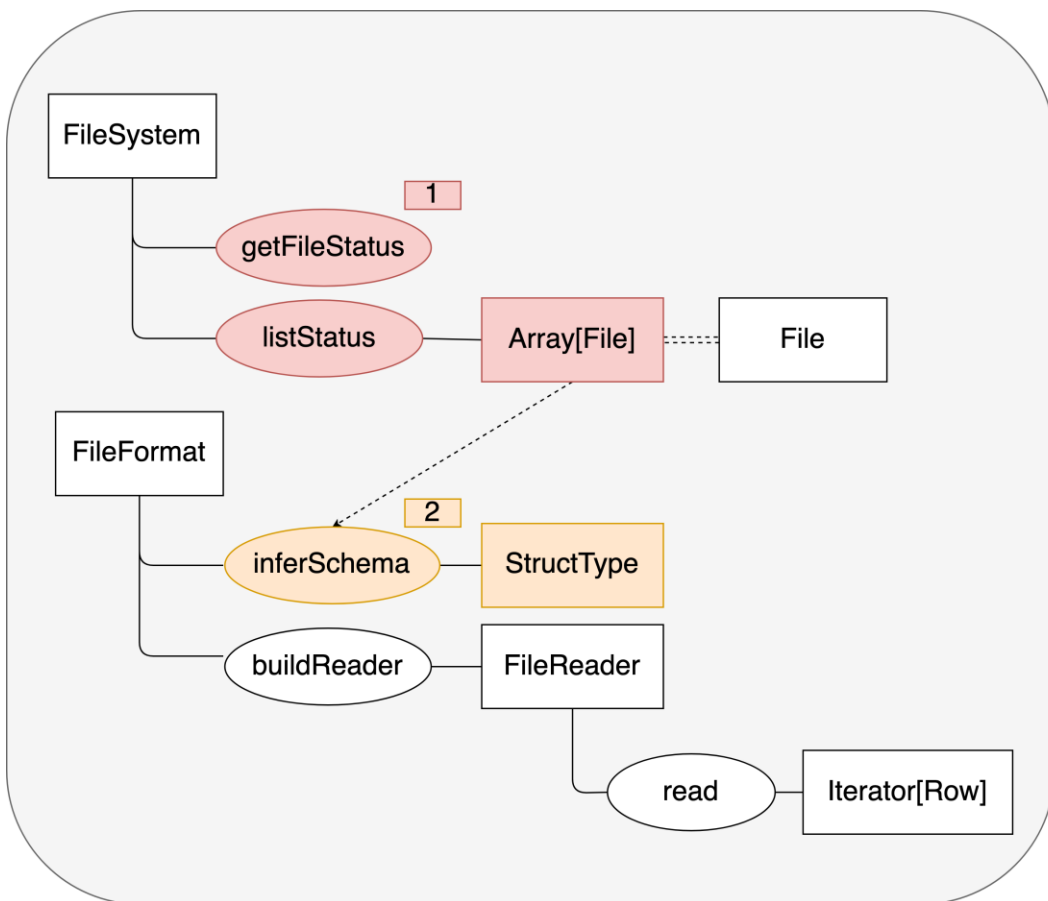
# Псевдолистинг таблицы



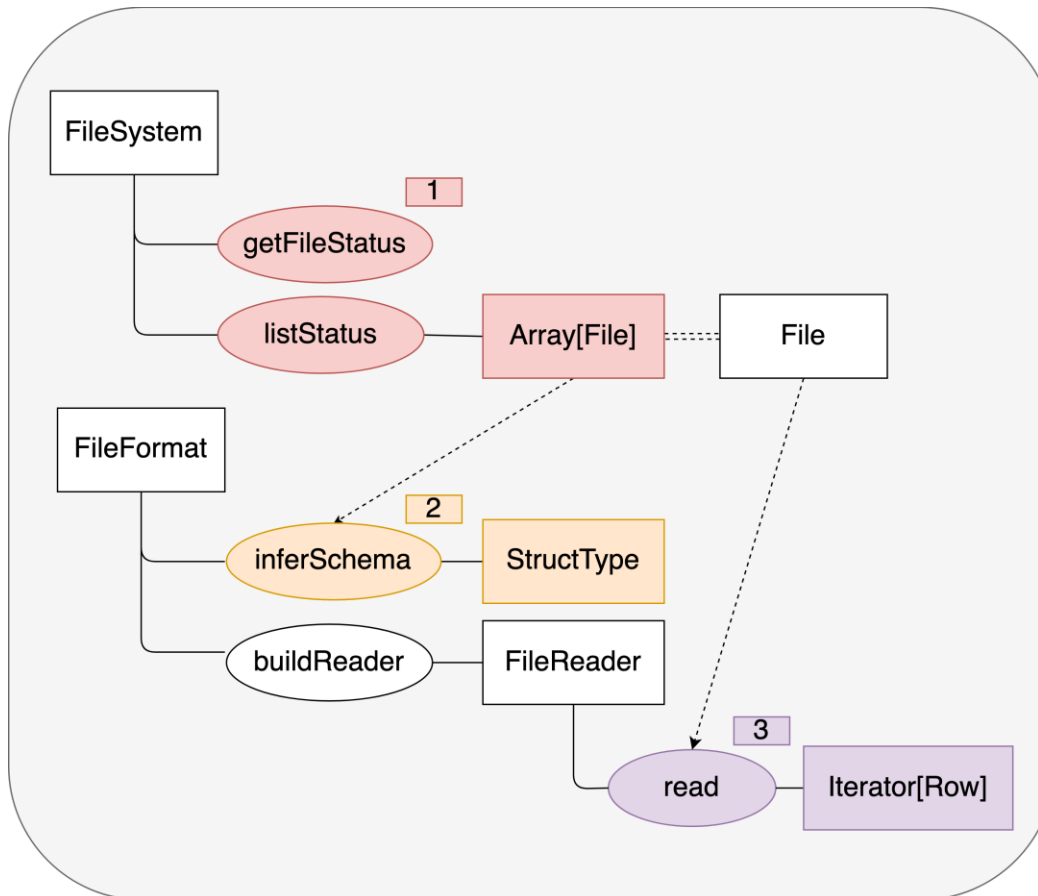
# Разбиваем таблицу на партии



# Читаем схему

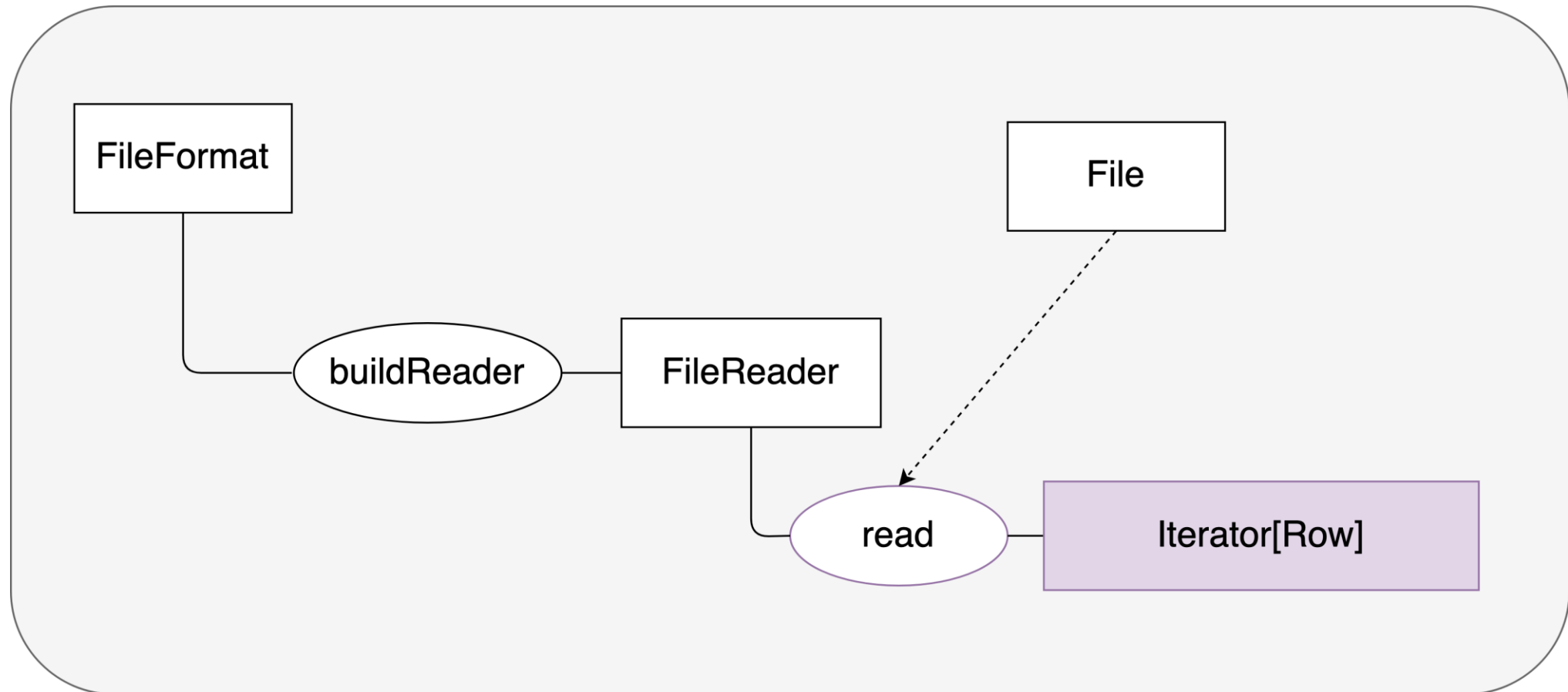


# Читаем отдельные партии

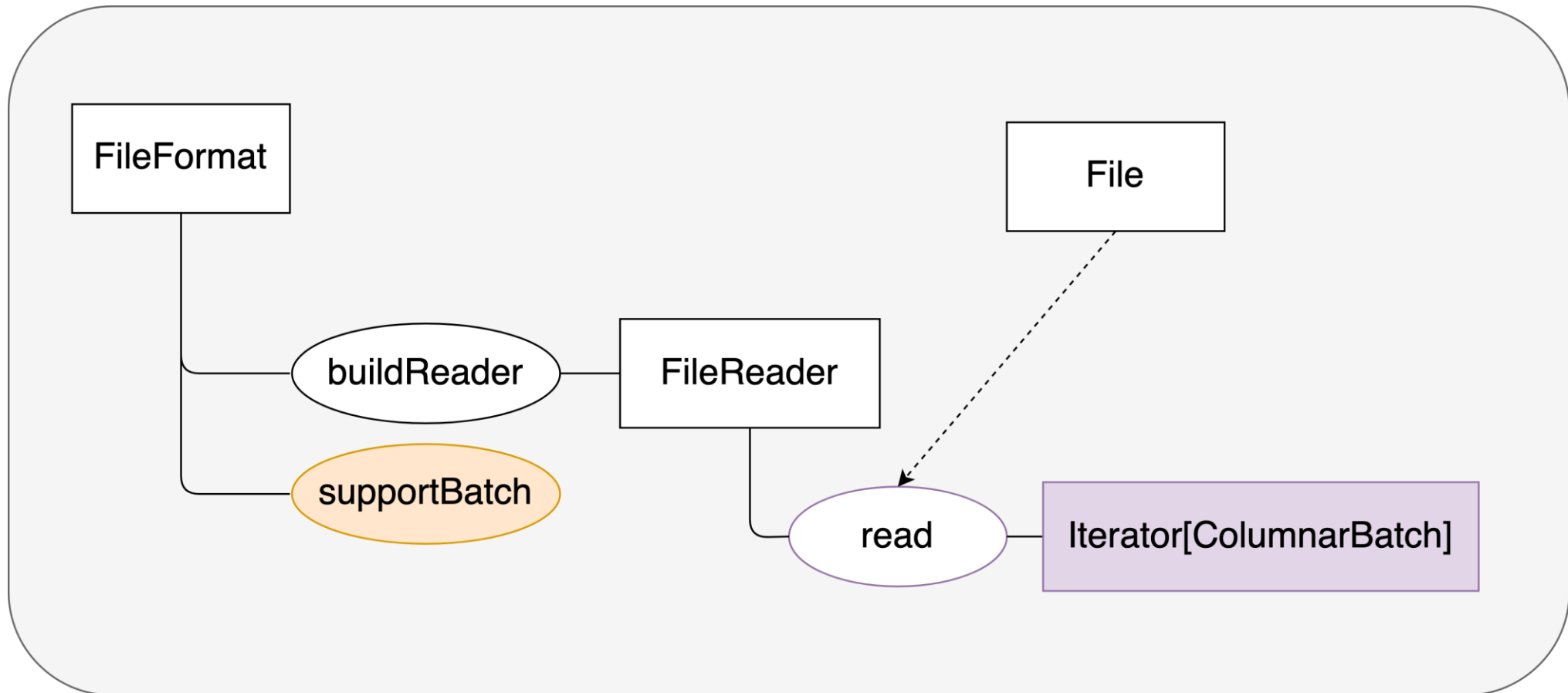


```
/
|---> tmp
|---> file.txt
|---> example
|----> table_1
|
|---> [0, 200]
|---> [201, 400]
|---> [401, 600]
|
|----> table_2
```

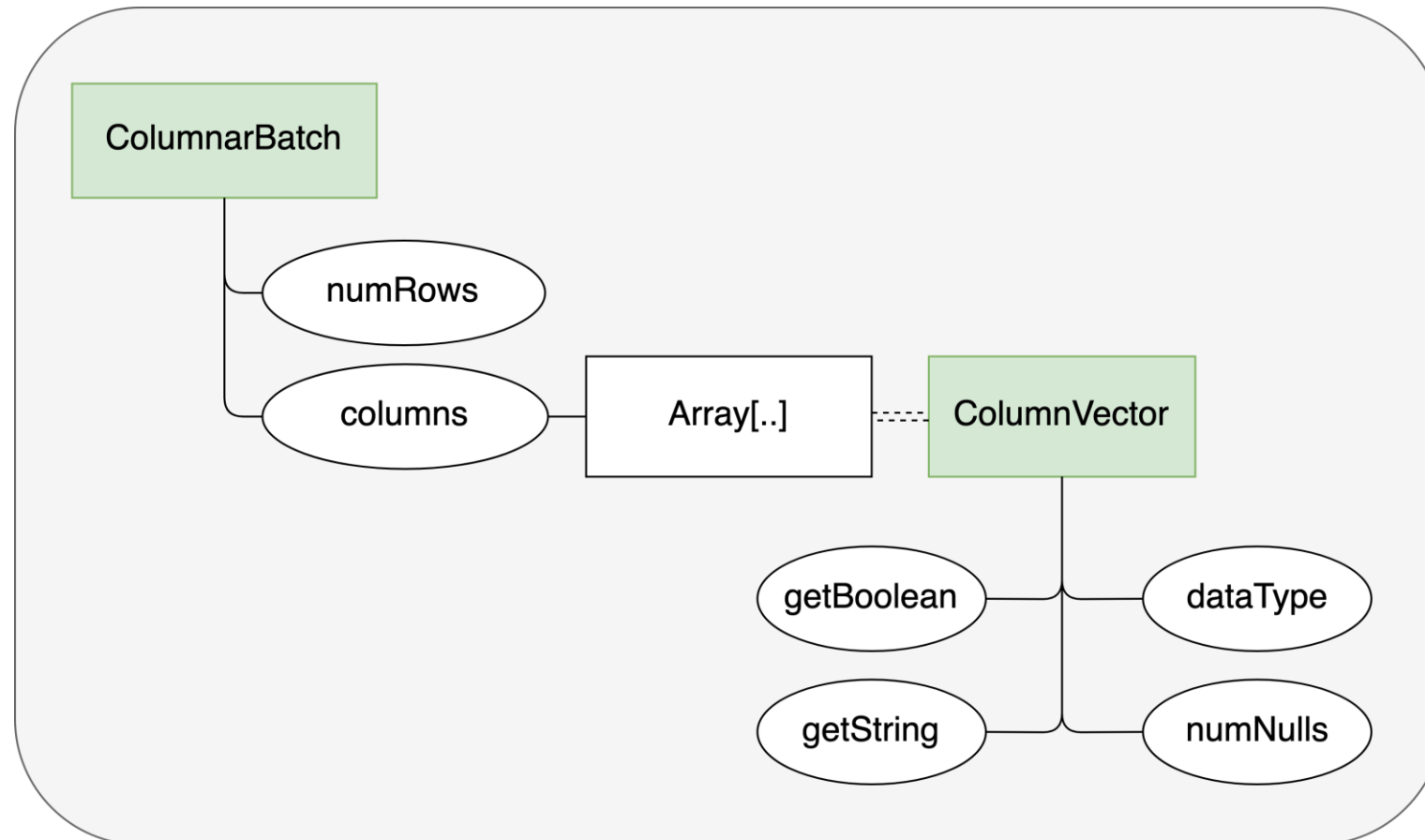
# Построчное чтение



# Чтение батчами




# ColumnarBatch





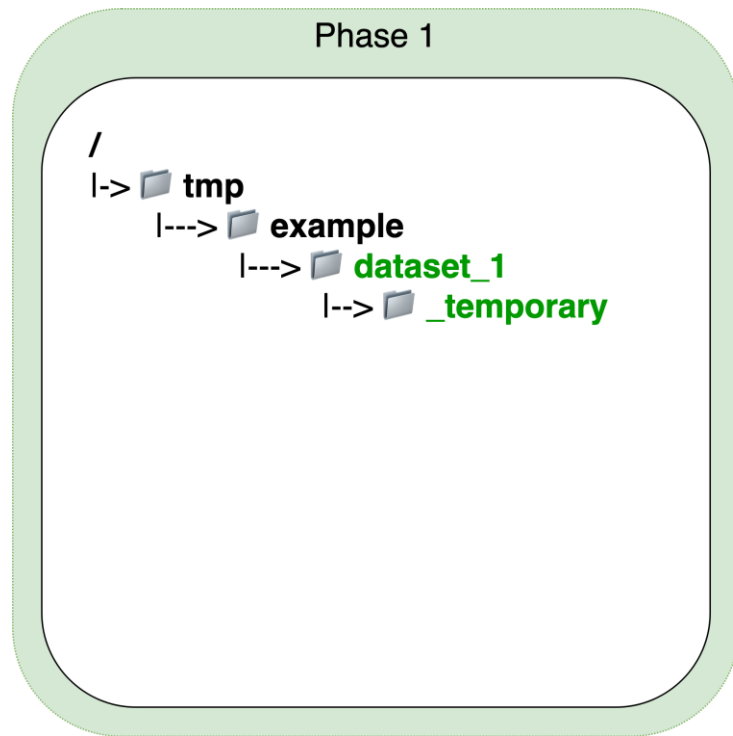
# Как Spark пишет Parquet в HDFS

# Запись в Parquet на HDFS

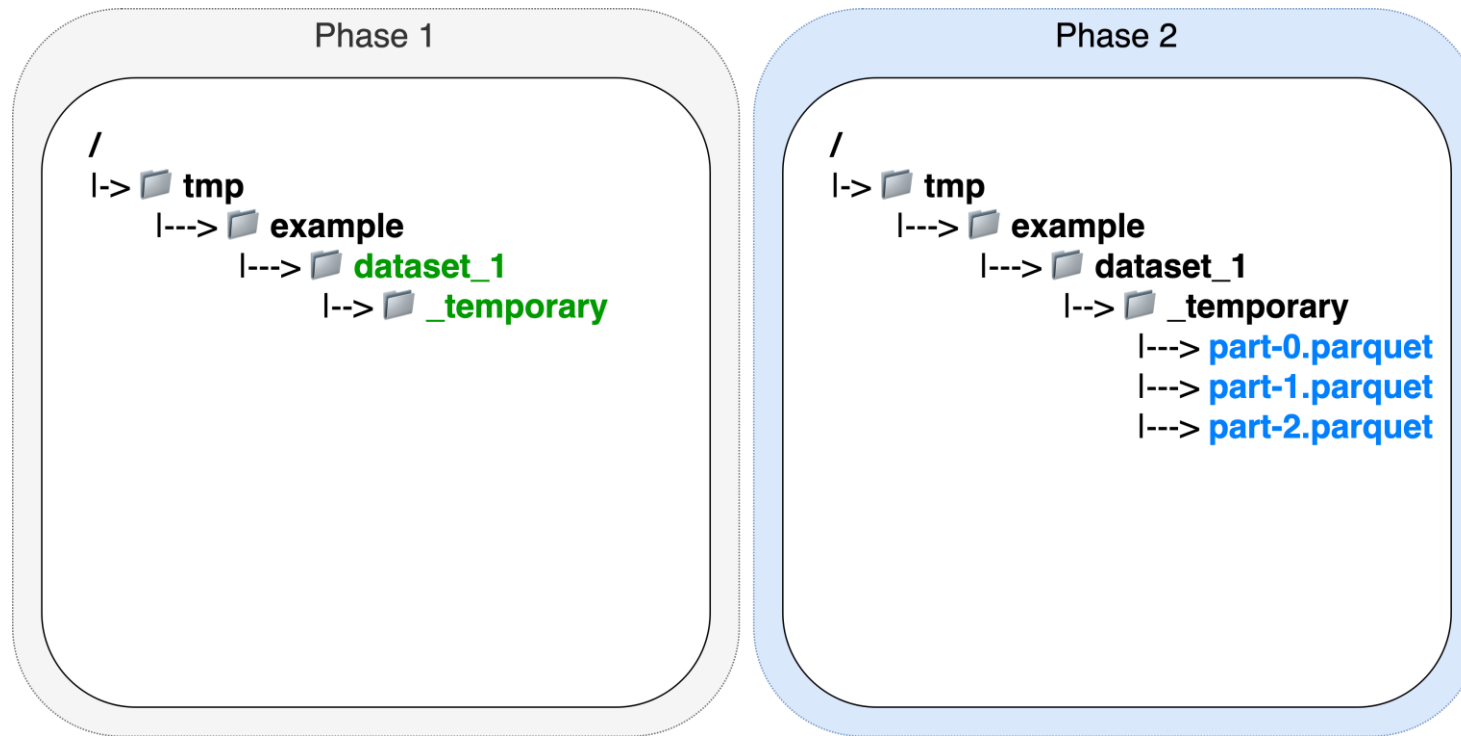
/  
|--->  tmp  
|--->  example

```
df.write.parquet("/tmp/example/dataset_1")
```

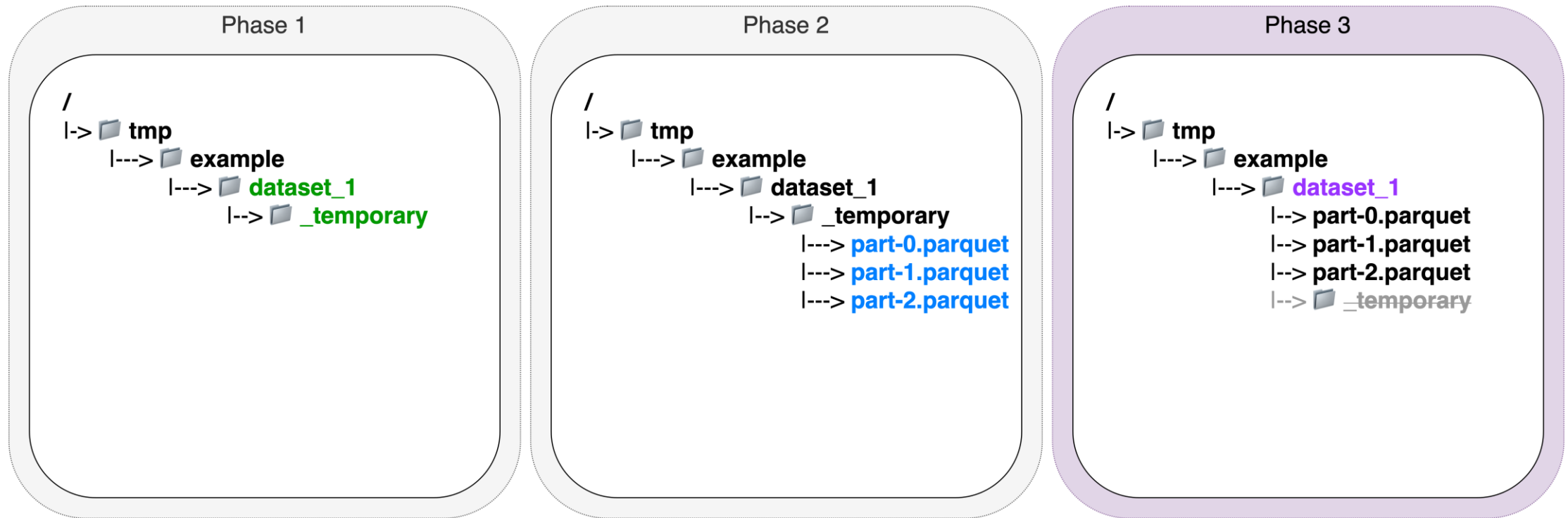
# Создаём временную директорию



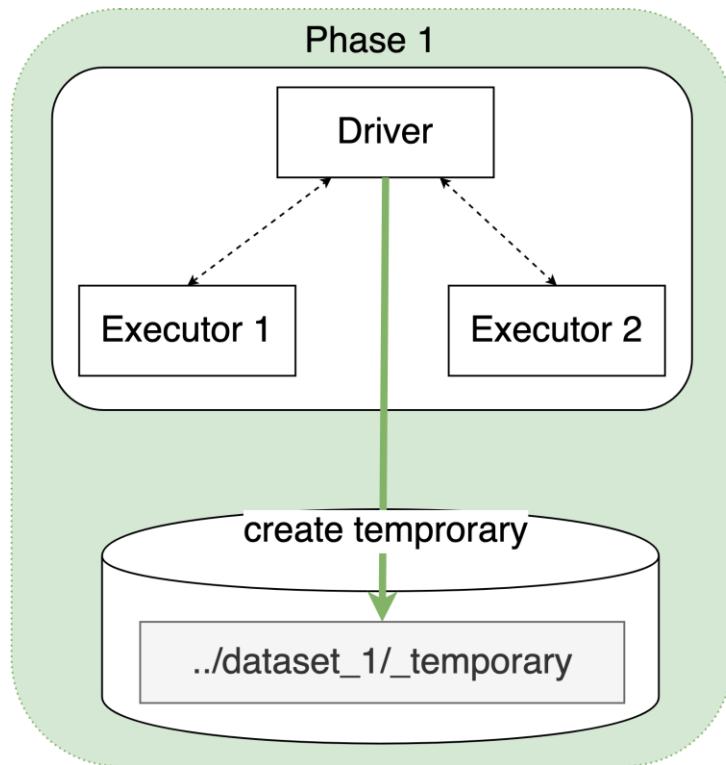
# Пишем файлы во временную директорию



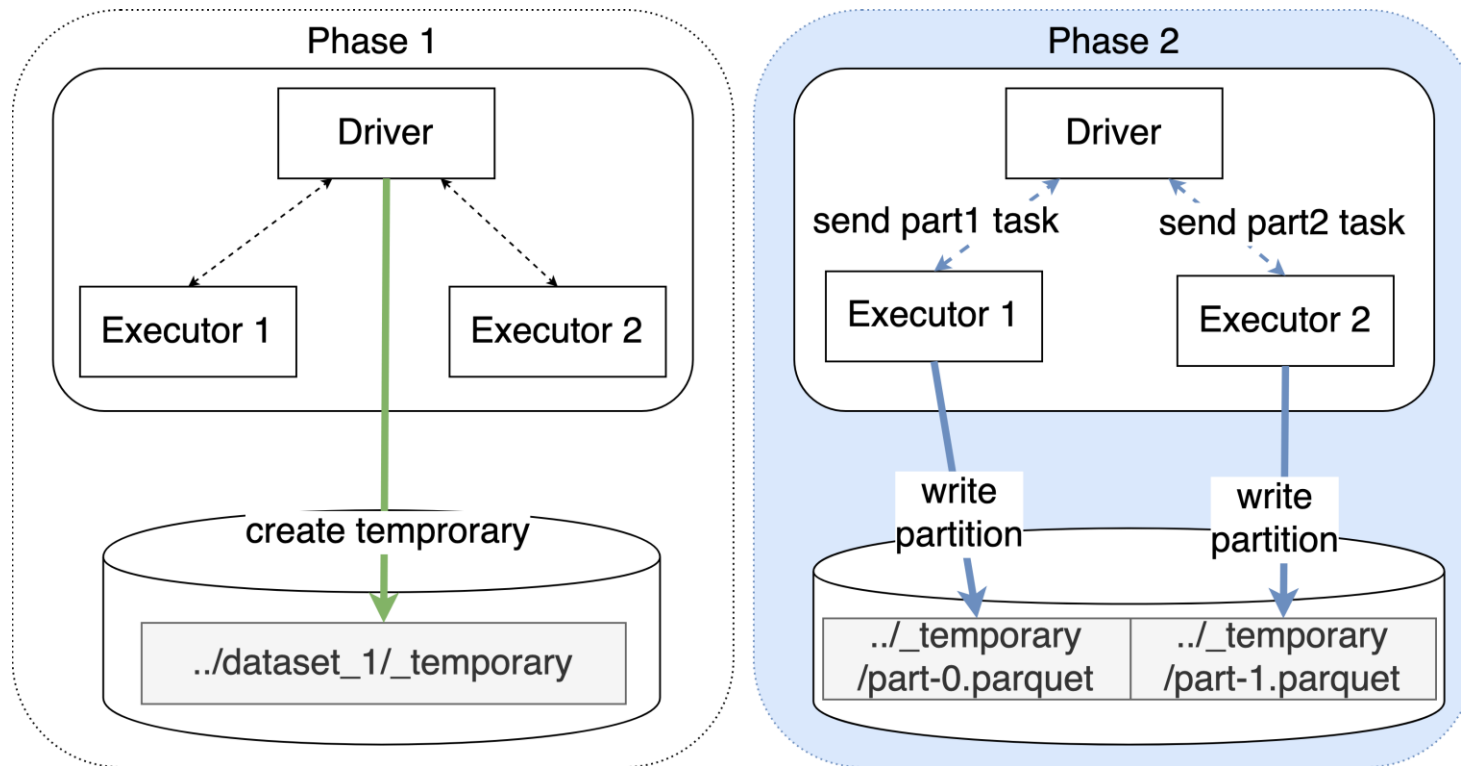
# Перемещаем файлы



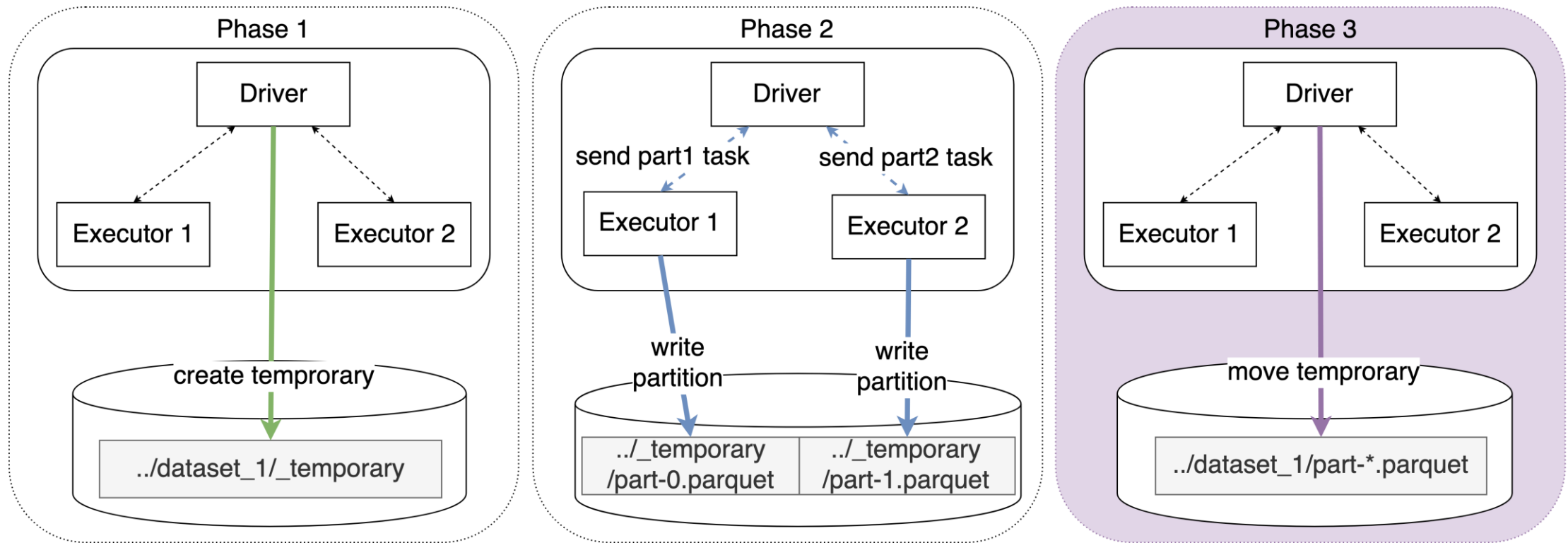
# Создаём временную директорию



# Пишем файлы во временную директорию

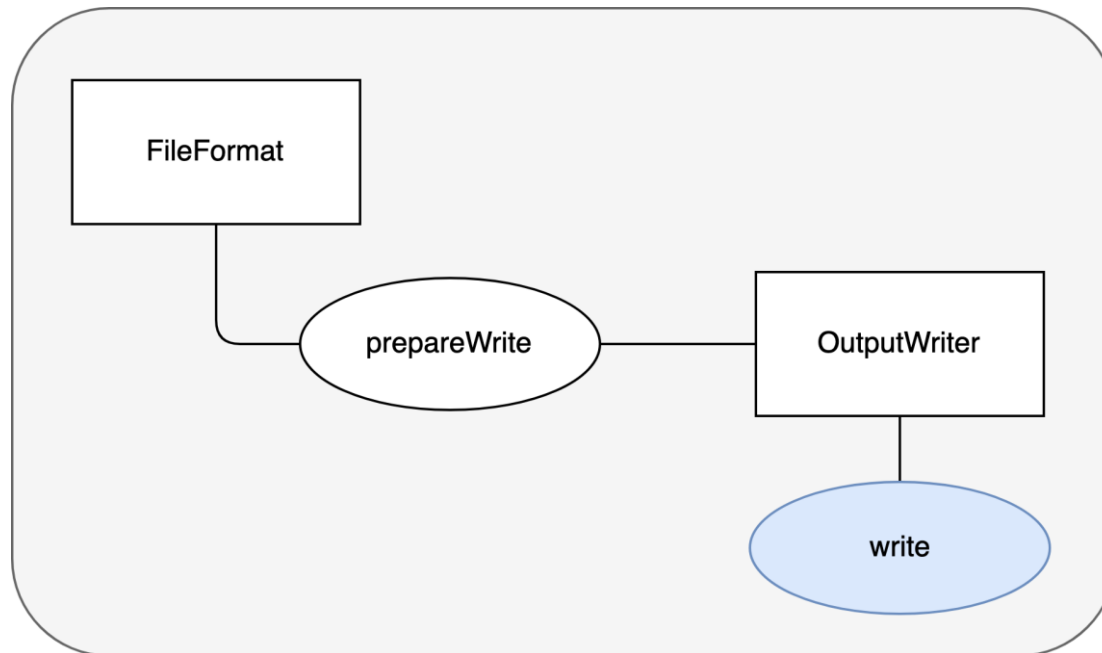


# Перемещаем файлы





# FileFormat

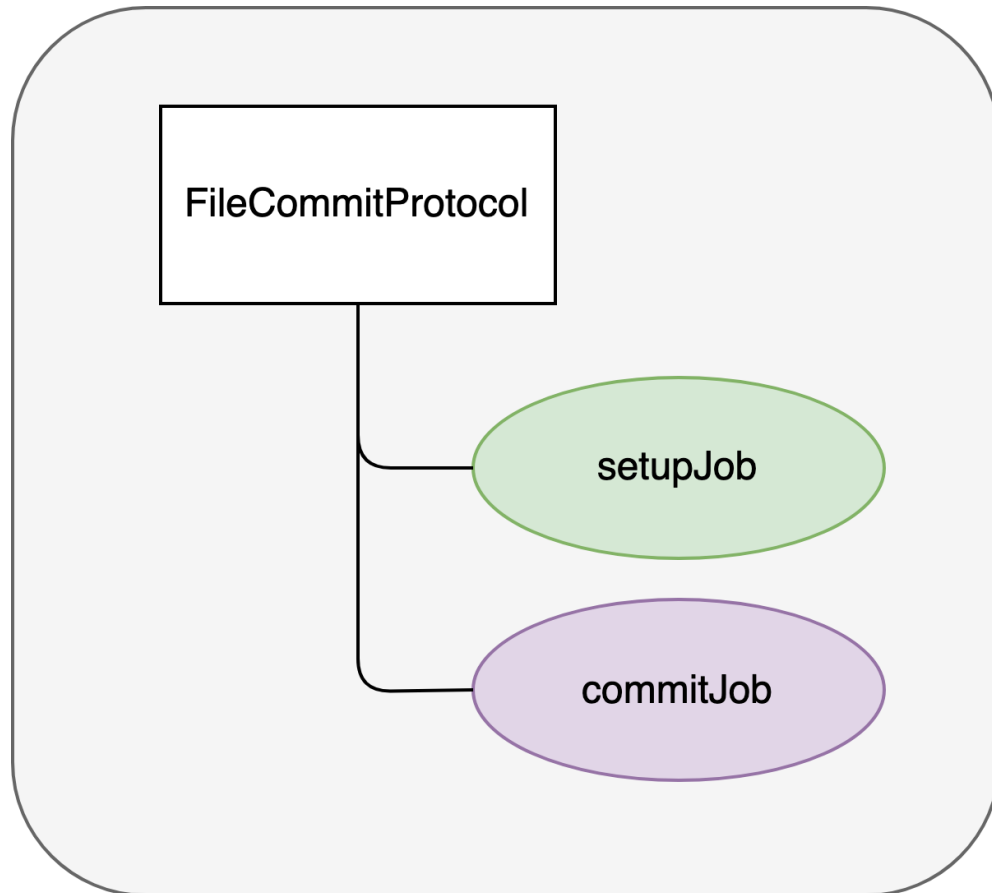


```
/
|-> tmp
|---> example
|---> dataset_1
|---> _temporary
```

```
/
|-> tmp
|---> example
|---> dataset_1
|---> _temporary
|---> part-0.parquet
|---> part-1.parquet
|---> part-2.parquet
```

```
/
|-> tmp
|---> example
|---> dataset_1
|---> part-0.parquet
|---> part-1.parquet
|---> part-2.parquet
|---> _temporary
```

# FileCommitProtocol

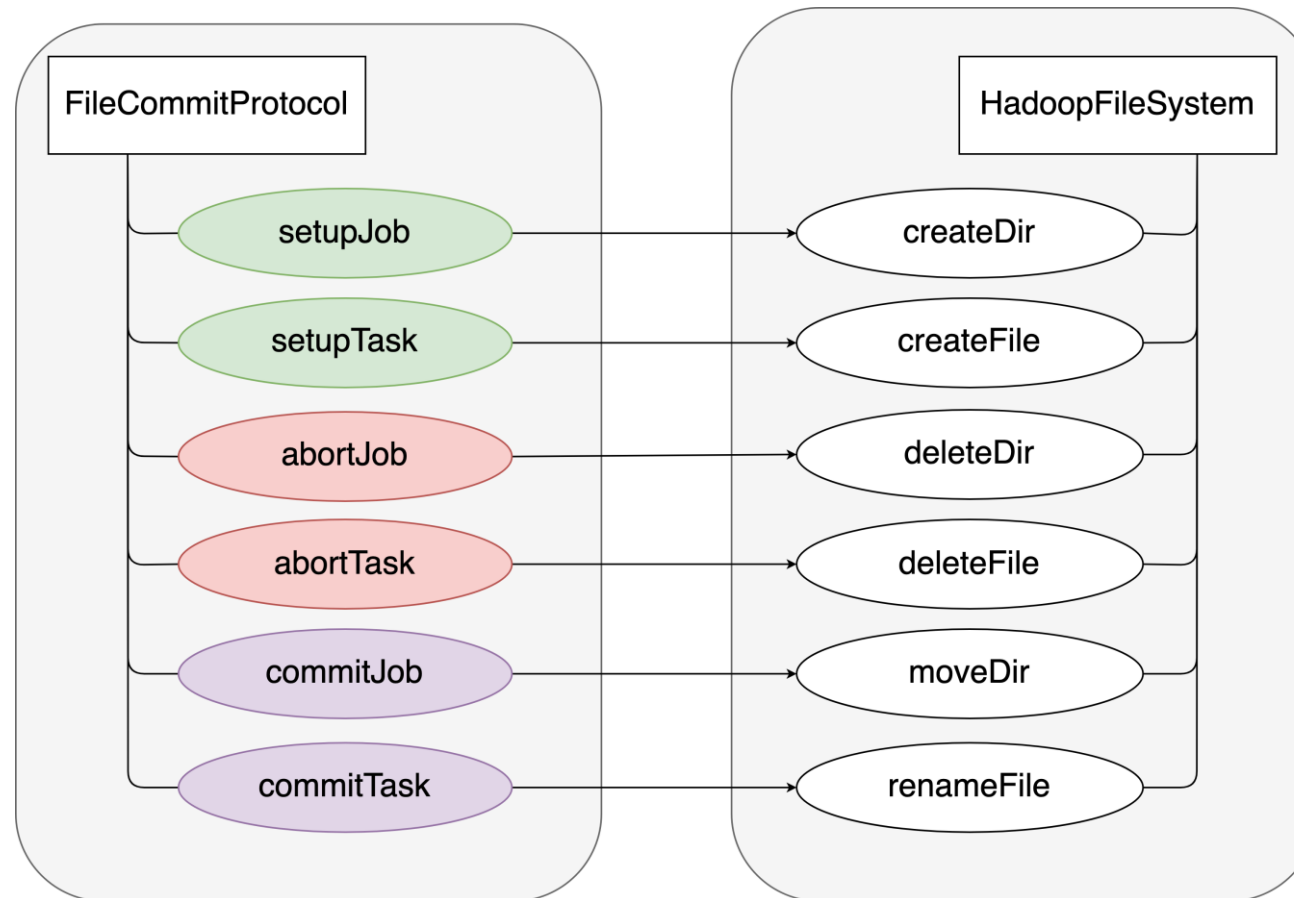


```
/
|-> tmp
|---> example
|---> dataset_1
|---> _temporary
```

```
/
|-> tmp
|---> example
|---> dataset_1
|---> _temporary
|---> part-0.parquet
|---> part-1.parquet
|---> part-2.parquet
```

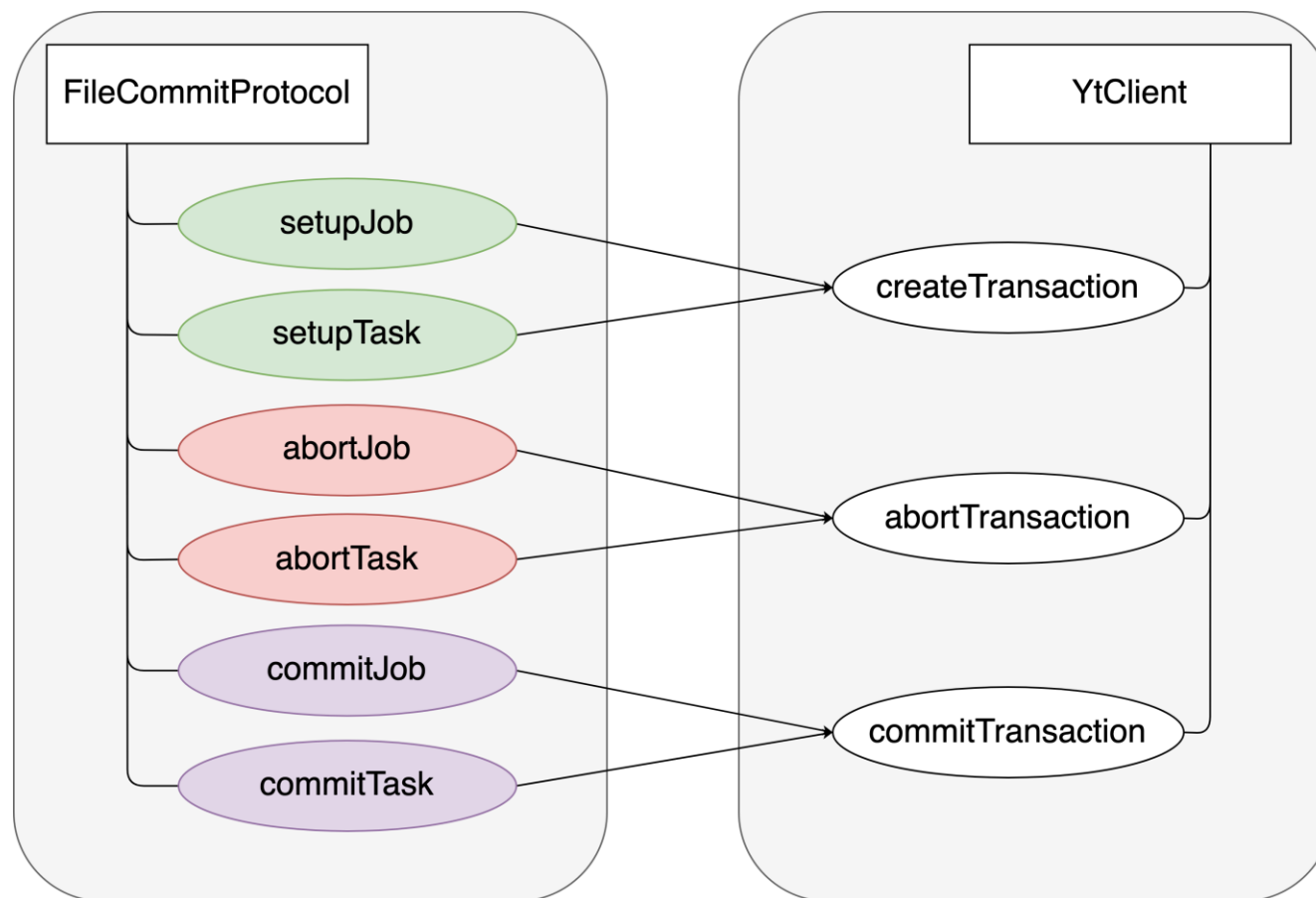
```
/
|-> tmp
|---> example
|---> dataset_1
|---> part-0.parquet
|---> part-1.parquet
|---> part-2.parquet
|---> _temporary
```

# Эмулируем транзакции в HDFS

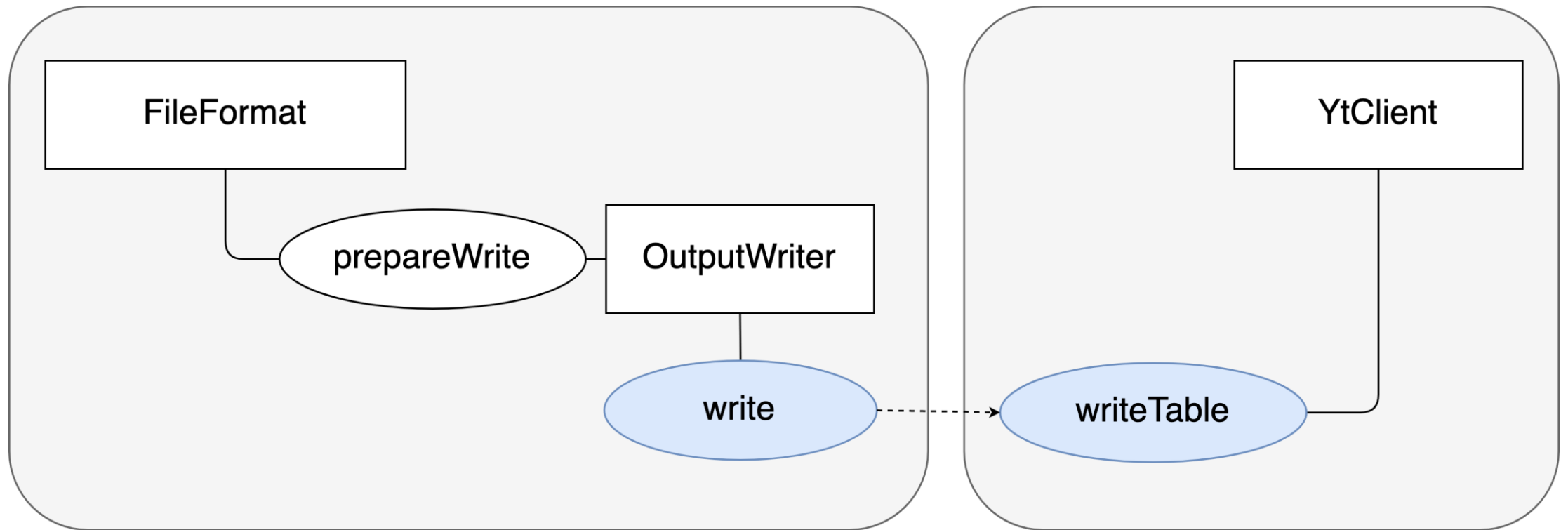


# Как SPYT пишет таблицы в YT



# В YТ уже есть транзакции!



# Запись партиций в YT



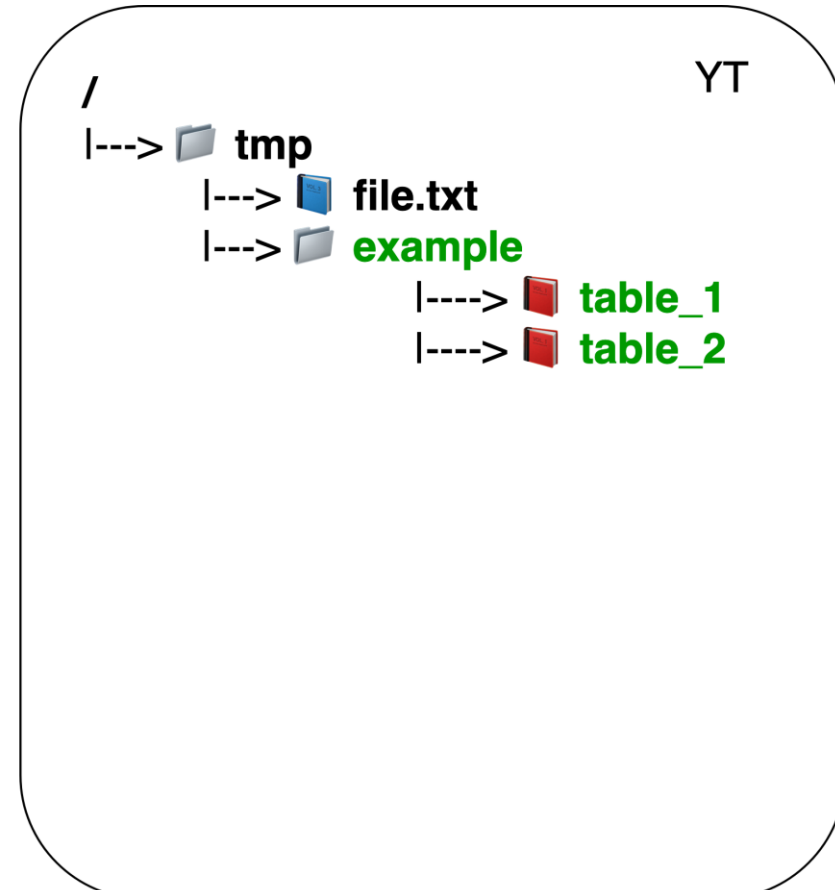
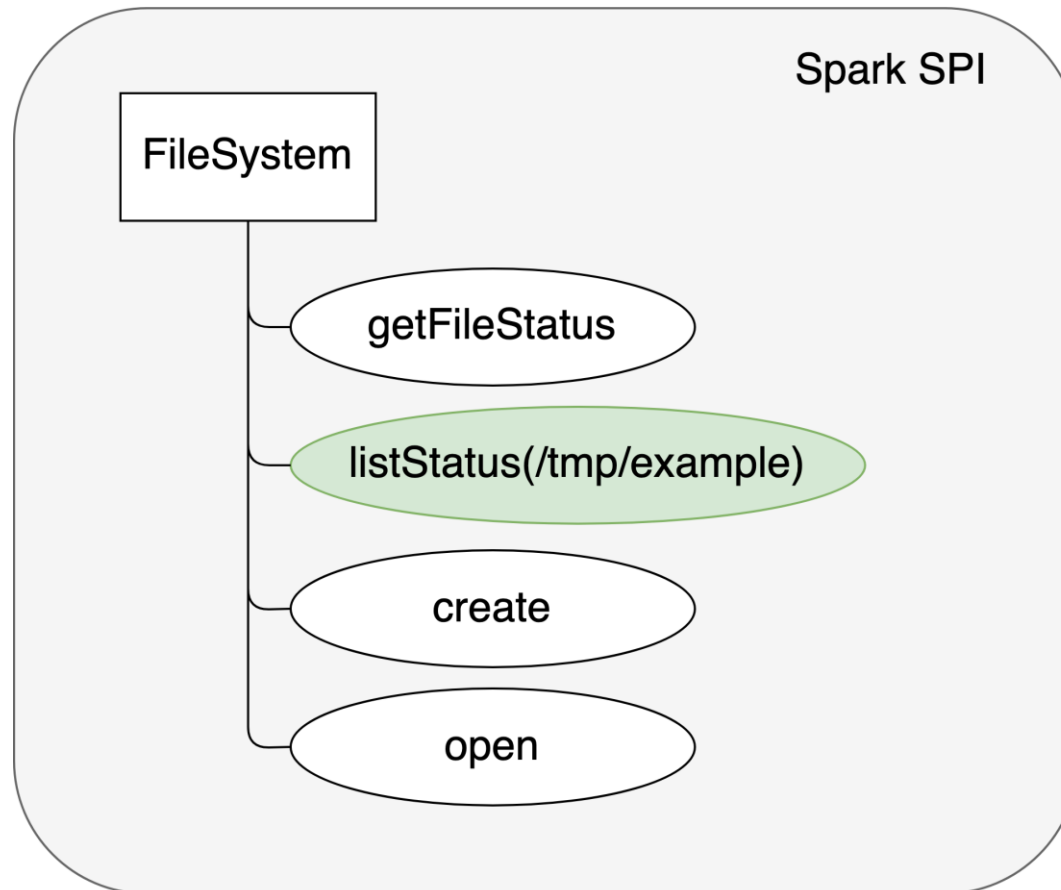
# Итоги продвинутого подхода

-  Научиться читать поколоночными батчами
-  Переиспользовать код, написанный для HDFS и Parquet:
  - Partition discovery
  - Обход директорий на чтении
  - Создание / удаление временных файлов на записи

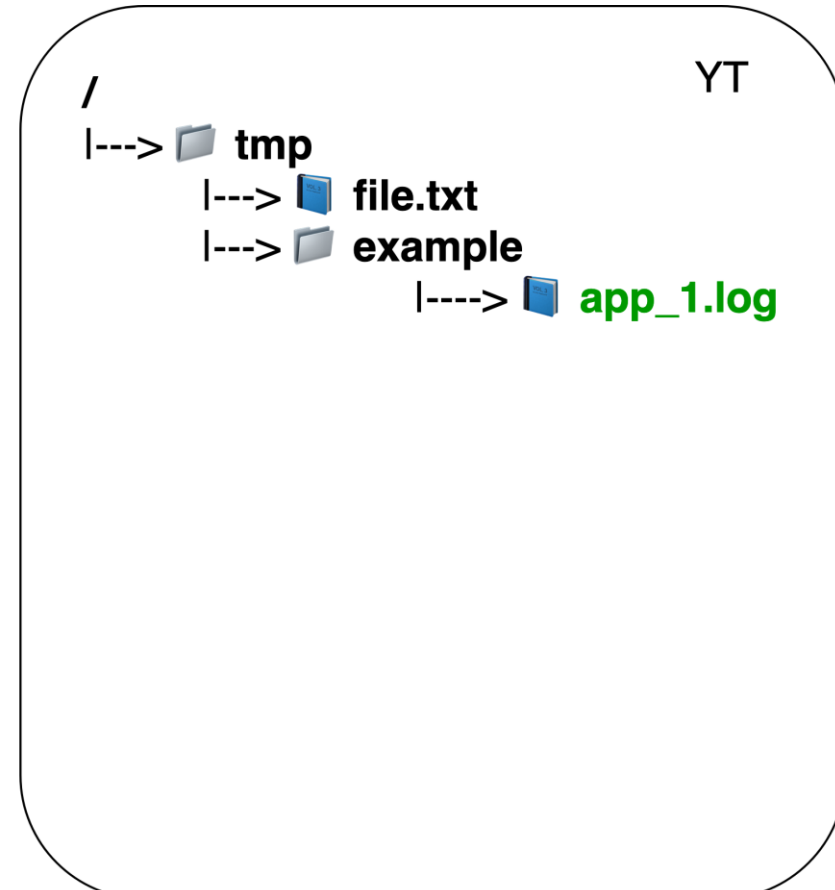
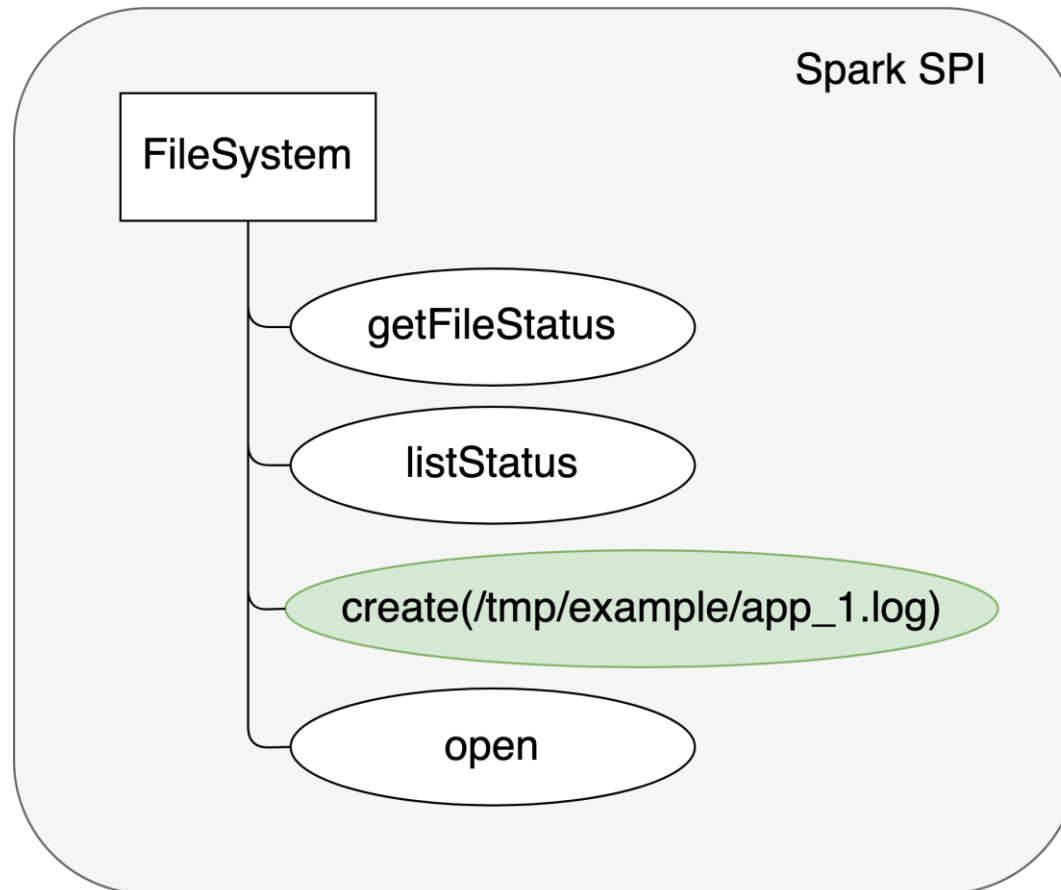
# Дополнительные фичи



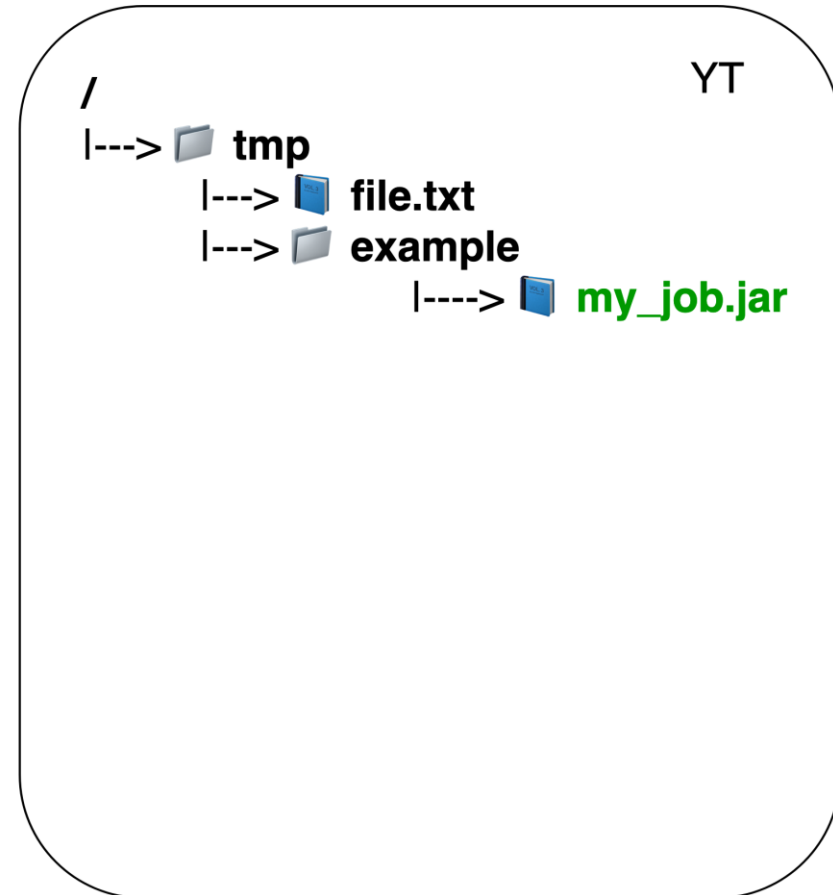
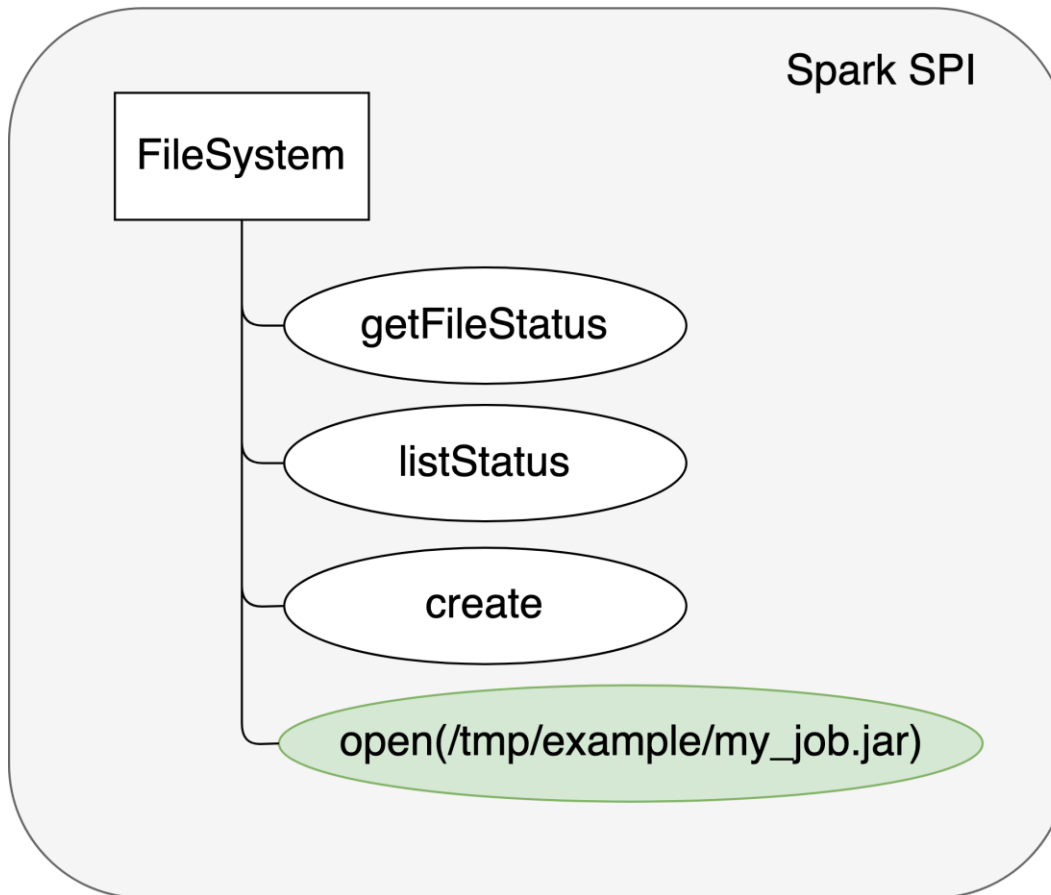
# FileSystem



# Создание файлов в YT



# Чтение файлов из YT

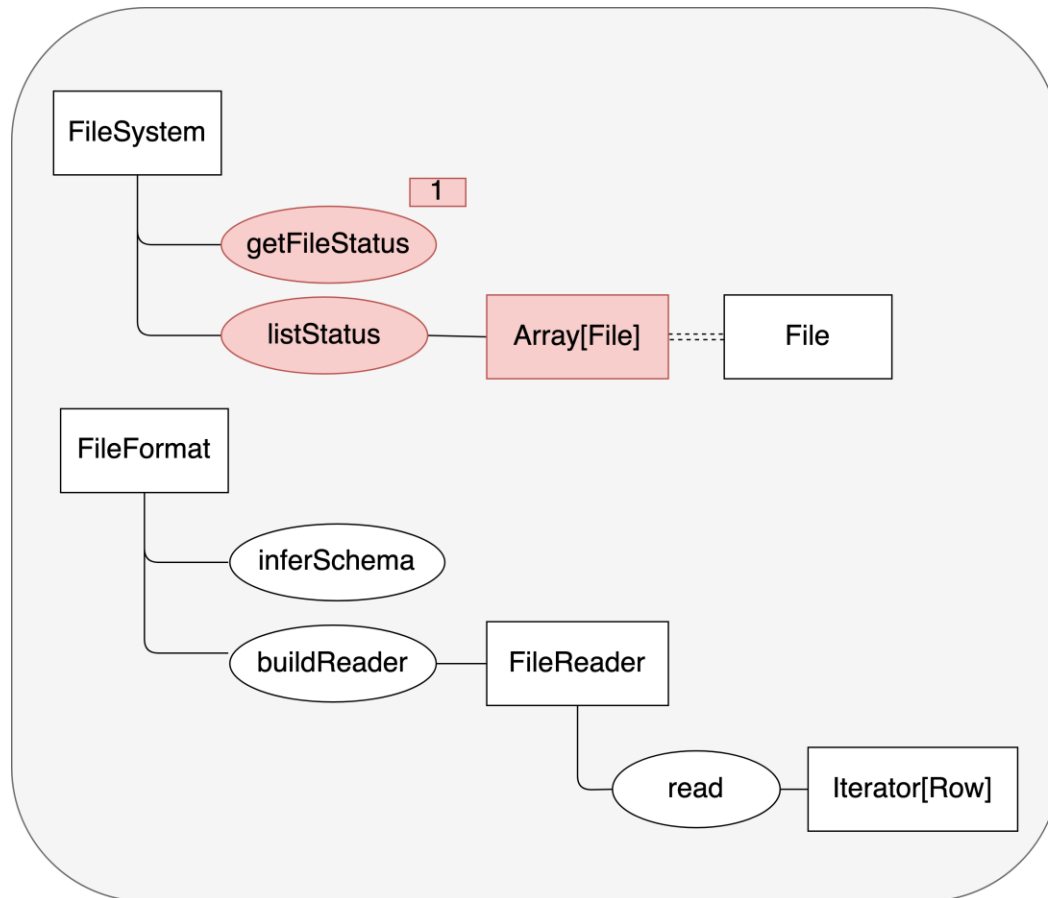


# Для чего это?

- Заработало скачивание файлов из YT в spark-submit
- Есть возможность сохранять event log в YT и читать его в SparkHistoryServer

# Грабли и костыли

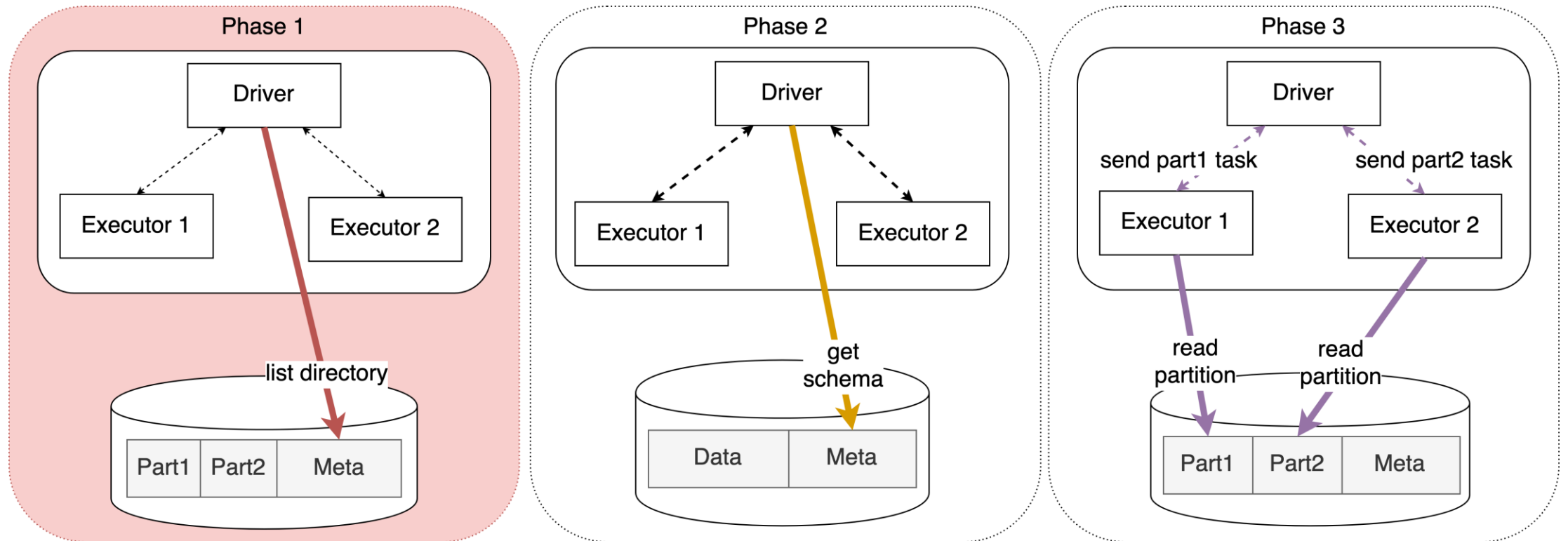
# Листинг таблицы при чтении



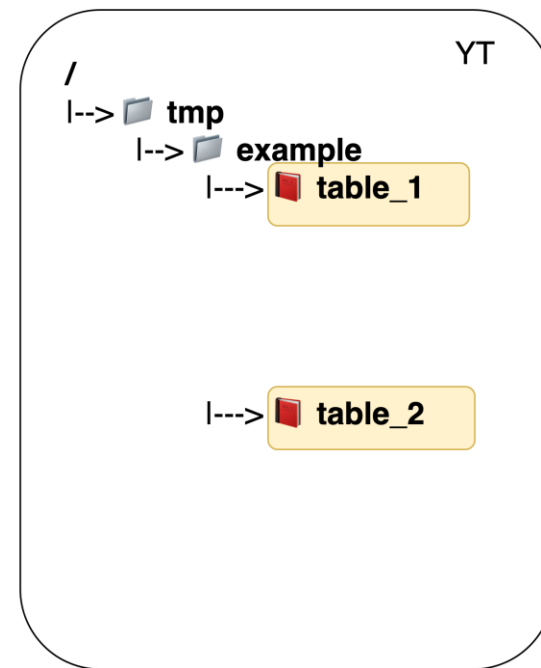
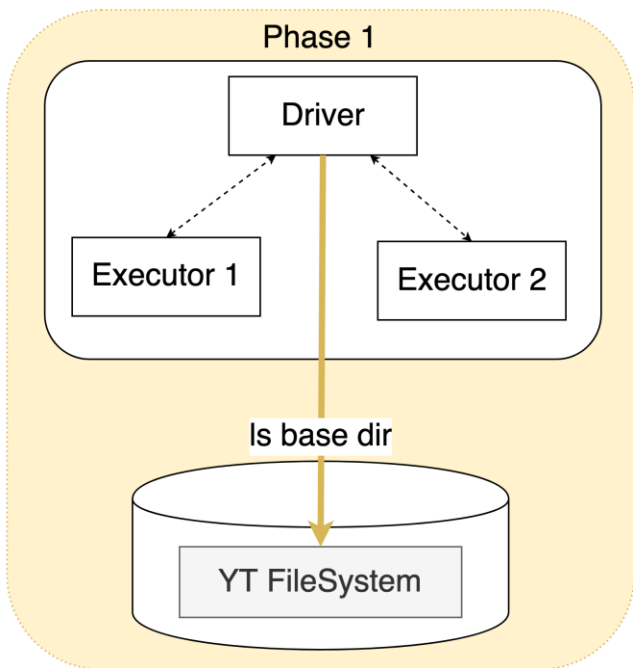
```
/
|---> tmp
|---> file.txt
|---> example
|---> table_1
|---> [0, 200]
|---> [201, 400]
|---> [401, 600]
|---> table_2
```

The diagram shows a tree structure of file paths. The root is **/**. It has three children: **tmp**, **file.txt**, and **example**. **example** has a child **table\_1**. **table\_1** has three children: **[0, 200]**, **[201, 400]**, and **[401, 600]**. **table\_2** is also a child of **/**.

# Листинг таблицы делает драйвер

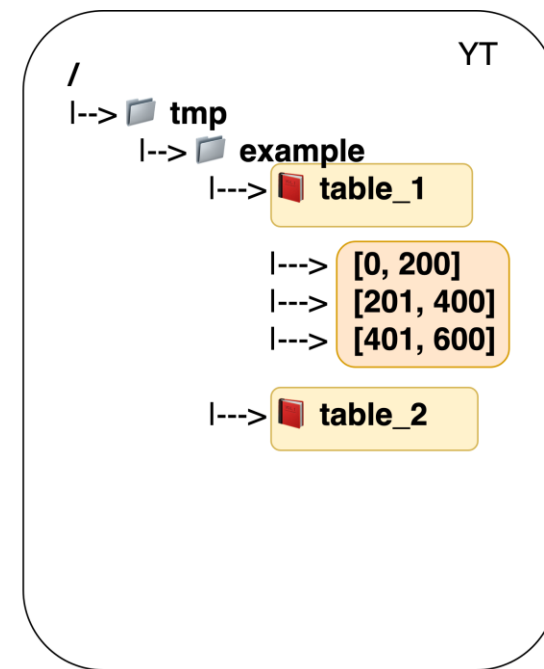
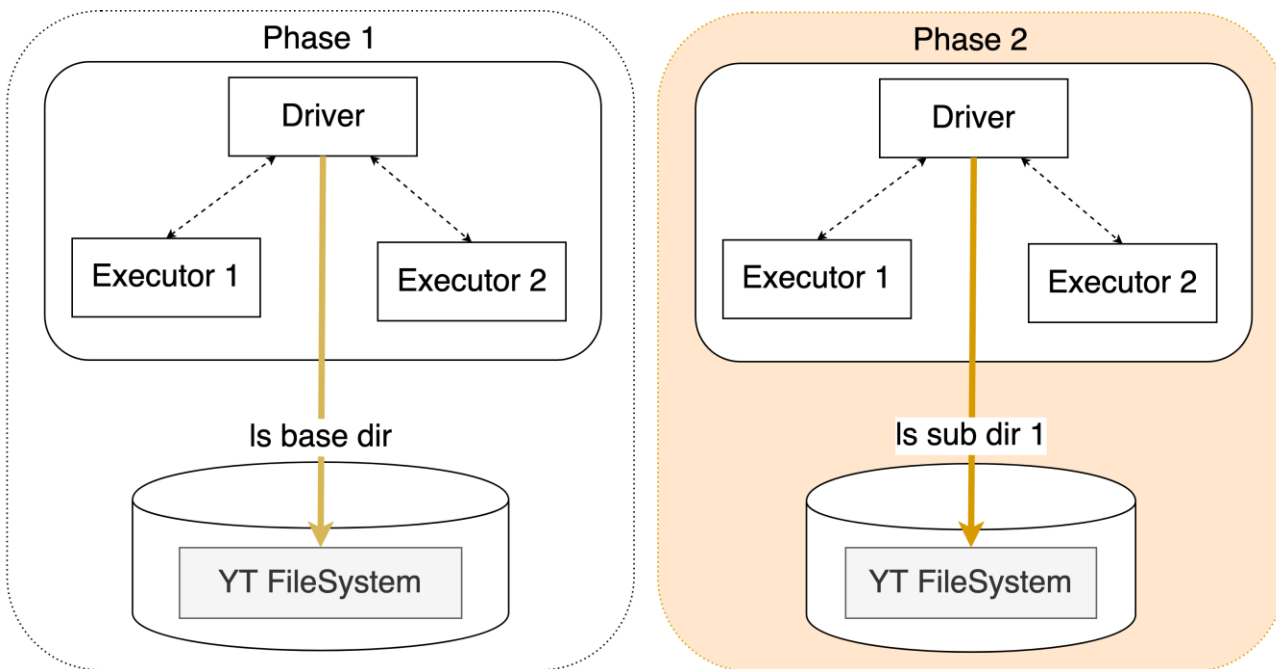


# Читаем /tmp/example целиком

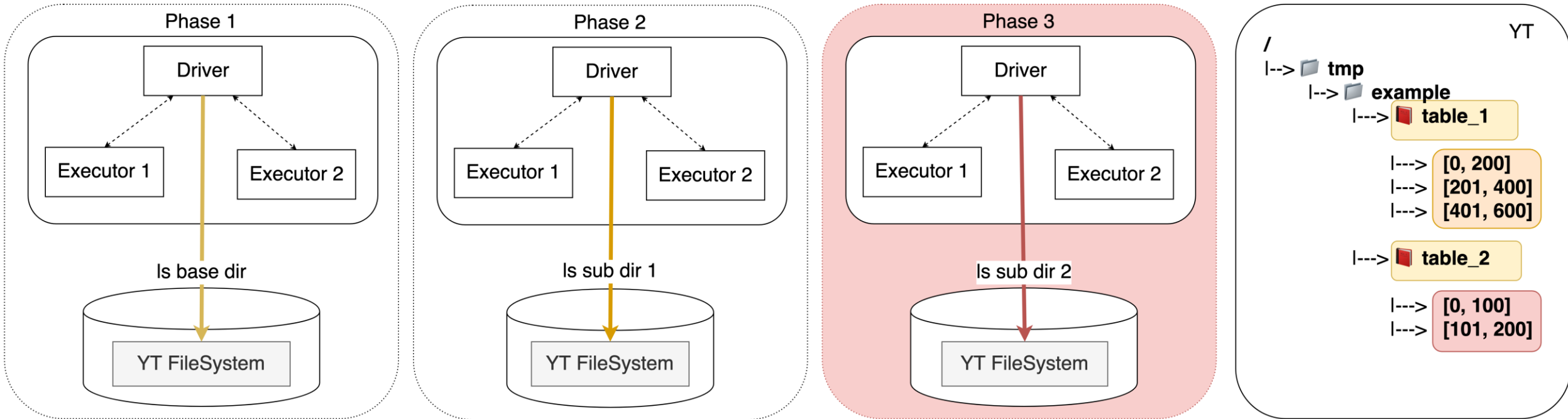




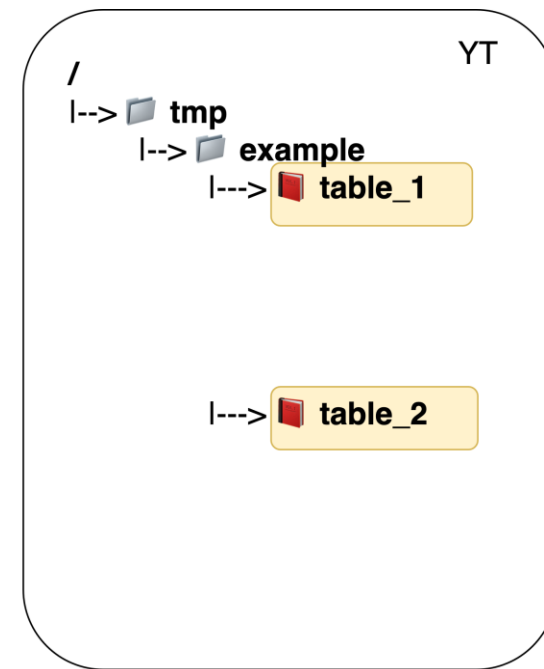
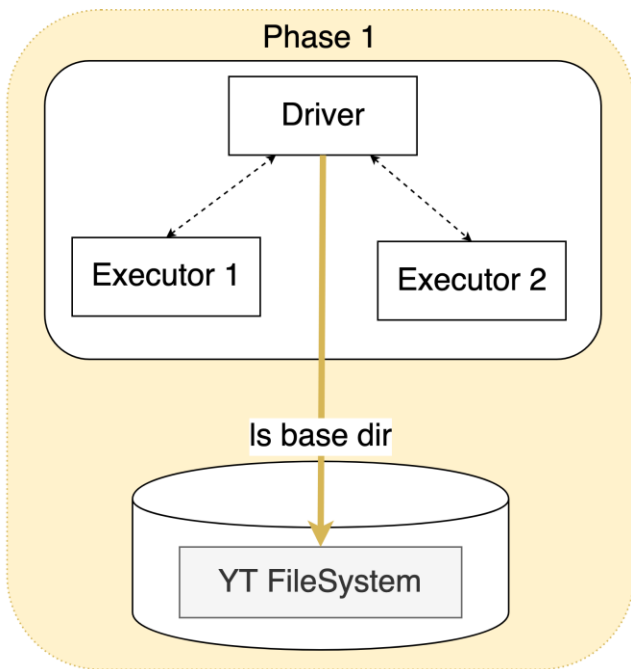
# Читаем /tmp/example целиком



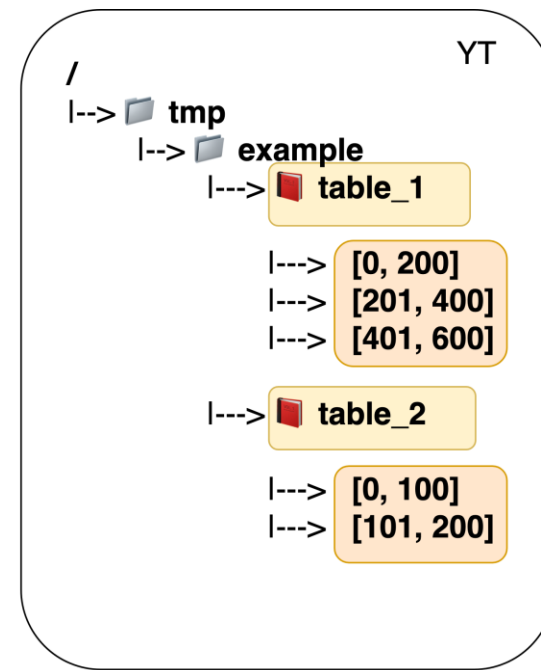
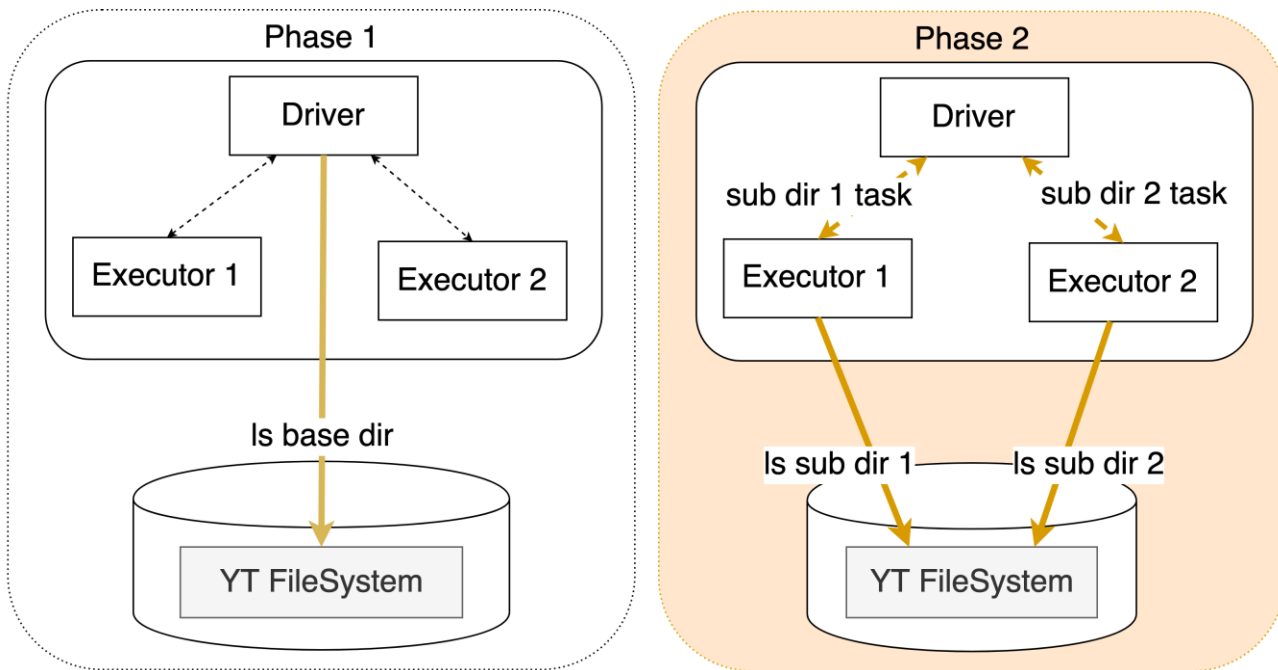
# Однопоточный листинг партиций



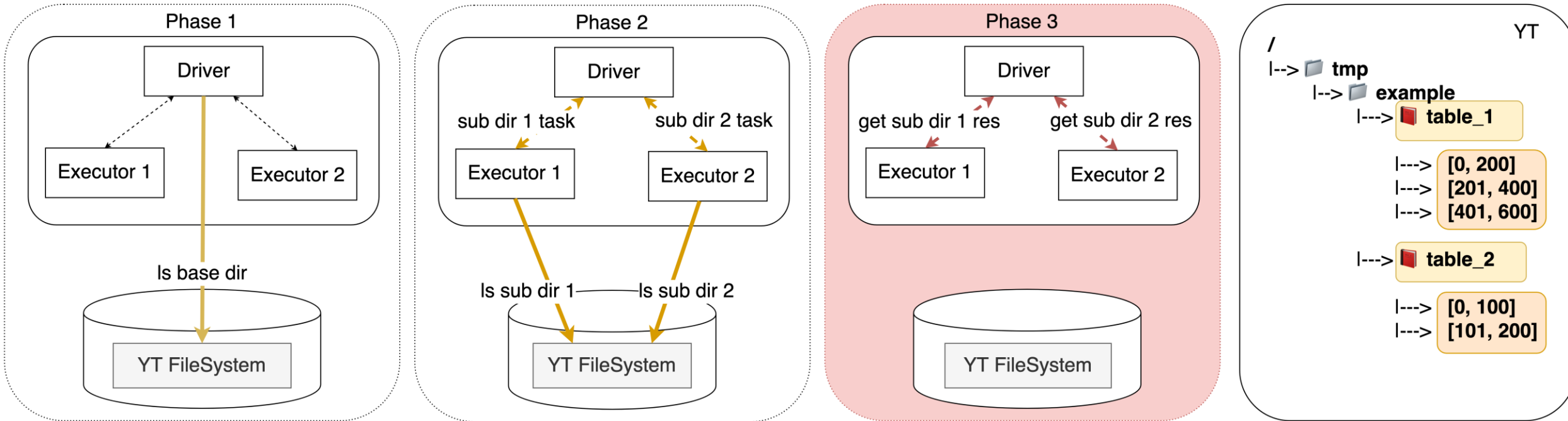
# Читаем /tmp/example целиком



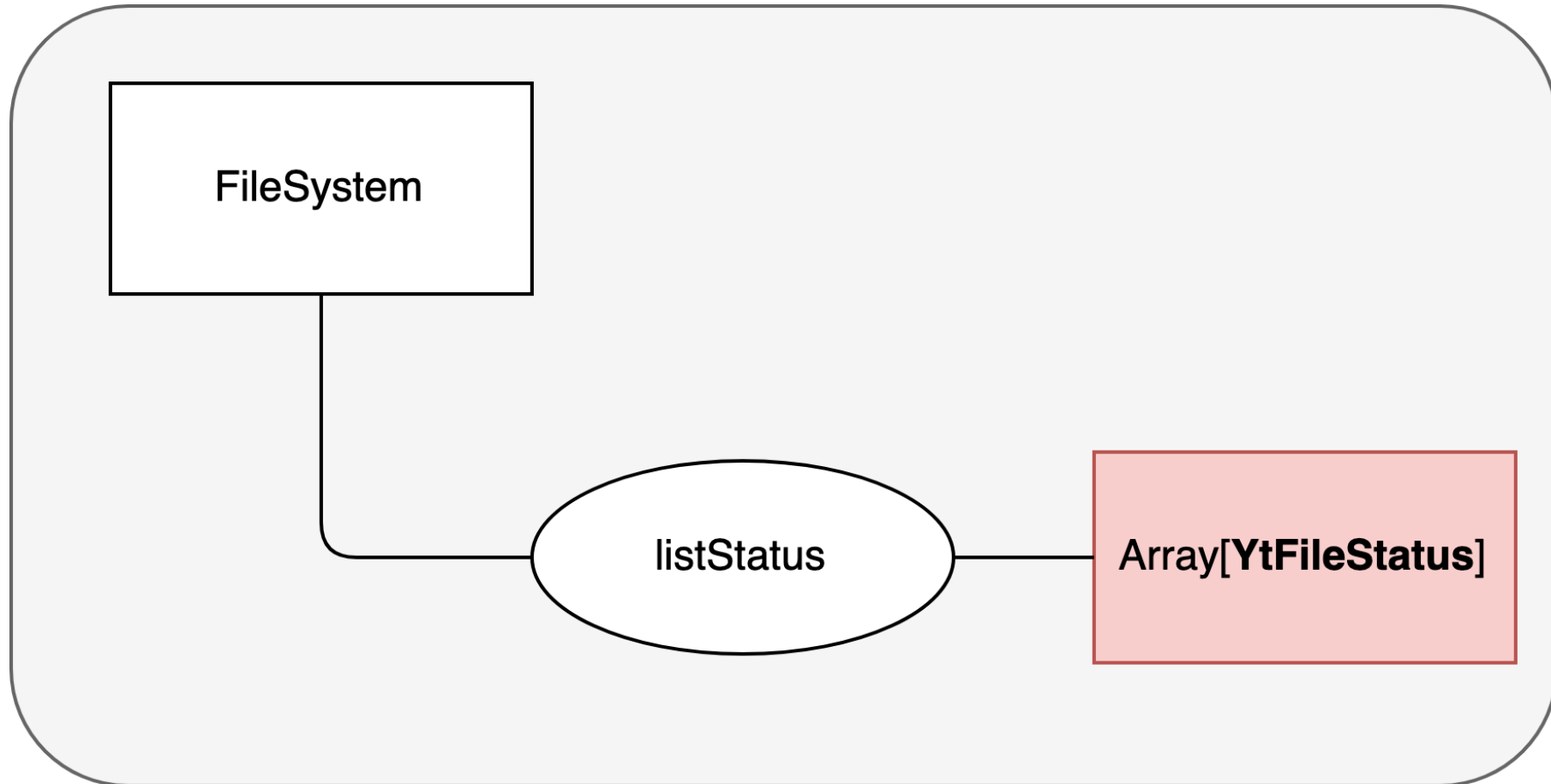
# Параллельный листинг на экзекьюторах



# Parallel partition discovery



# Реализация FileSystem



# Что с этим не так?

listStatus вызывается на экзекьюторе

# Что с этим не так?

listStatus вызывается на экзекьюторе



На экзекьюторе получаем Array[YtFileStatus]



# Что с этим не так?

listStatus вызывается на экзекьюторе



На экзекьюторе получаем Array[YtFileStatus]



Сериализуем и отправляем на драйвер

# Что с этим не так?

listStatus вызывается на экзекьюторе



На экзекьюторе получаем Array[YtFileStatus]



Сериализуем и отправляем на драйвер



Десериализуем на драйвере

# Что с этим не так?

listStatus вызывается на экзекьюторе



На экзекьюторе получаем Array[YtFileStatus]



Сериализуем и отправляем на драйвер



Десериализуем на драйвере



Array[FileStatus]

# Что делать?

- Учитывать, что методы для листинга директории могут вызываться на экзекьюторах
- Не пытаться обмануть SPI и использовать FileStatus

# FileSystem – синглтон?

```
public static FileSystem get (URI uri,  
                             Configuration conf)  
{  
    ...  
    return CACHE.get (uri, conf);  
}
```

# Использую синглтон в init / close

```
class YtFileSystem extends FileSystem {  
  override def initialize(...): Unit = {  
    ...  
    initYtClient()  
  }  
  override def close(): Unit = {  
    ...  
    closeYtClient()  
  }  
  ...  
}
```

# Как устроен кэш FileSystem

```
private FileSystem getInternal ( Key key, ...) {  
    ...  
    synchronized (this) {  
        fs = map.get(key);  
    }  
    if (fs != null) return fs;  
    fs = createFileSystem(uri, conf);  
  
    synchronized (this) {  
        // refetch the lock again  
        FileSystem oldfs = map.get(key);  
        if (oldfs != null) {  
            // FS created while lock is releasing  
            fs.close();  
            // close the new file system  
            return oldfs;  
            // return the old file system  
        }  
    }  
    ...  
    return fs;  
}  
}
```

# Пытаемся взять FS из кэша

```
private FileSystem getInternal ( Key key, ...) {  
    ...  
    synchronized (this) {  
        fs = map.get(key);  
    }  
    if (fs != null) return fs;  
    fs = createFileSystem(uri, conf);  
  
    synchronized (this) {  
        // refetch the lock again  
        FileSystem oldfs = map.get(key);  
        if (oldfs != null) {  
            // FS created while lock is releasing  
            fs.close();  
            // close the new file system  
            return oldfs;  
            // return the old file system  
        }  
    }  
    ...  
    return fs;  
}  
}
```



# Создаём новый инстанс FS

```
private FileSystem getInternal ( Key key, ...) {
    ...
    synchronized (this) {
        fs = map.get(key);
    }
    if (fs != null) return fs;

    fs = createFileSystem(uri, conf);

    synchronized (this) { // refetch the lock again
        FileSystem oldfs = map.get(key);
        if (oldfs != null) { // FS created while lock is releasing
            fs.close(); // close the new file system
            return oldfs; // return the old file system
        }
    }

    ...
    return fs;
}
}
```

# Другой поток создал инстанс раньше

```
private FileSystem getInternal ( Key key, ...) {
    ...
    synchronized (this) {
        fs = map.get(key);
    }
    if (fs != null) return fs;

    fs = createFileSystem(uri, conf);

    synchronized (this) { // refetch the lock again
        FileSystem oldfs = map.get(key);
        if (oldfs != null) { // FS created while lock is releasing
            fs.close(); // close the new file system
            return oldfs; // return the old file system
        }
    }

    ...
    return fs;
}
}
```

# Закрываем свой инстанс

```
private FileSystem getInternal ( Key key, ...) {  
    ...  
    synchronized (this) {  
        fs = map.get(key);  
    }  
    if (fs != null) return fs;  
    fs = createFileSystem(uri, conf);  
  
    synchronized (this) {  
        // refetch the lock again  
        FileSystem oldfs = map.get(key);  
        if (oldfs != null) {  
            // FS created while lock is releasing  
            fs.close();  
            // close the new file system  
            return oldfs;  
            // return the old file system  
        }  
    }  
    ...  
    return fs;  
}  
}
```

# Что делать?

- Запомнить, что FileSystem может быть создан не один, а лишний сразу будет закрыт
- Не использовать настоящие синглтоны в коде FileSystem

# Кастомные параметры в Spark

- Большую часть параметров подключения можно прокидывать через конфиг спарка:

```
--conf spark.my.param=value
```

- Их удобно доставать из SparkConf, можно определять дефолтные значения в `spark-defaults.conf`
- Но они видны в логах и на вкладке Environment в Spark UI

# Как скрыть отображение секретных параметров?

- В Spark есть специальный конфиг

```
spark.redaction.regex
```

- Все параметры, названия которых соответствуют регулярке, будут скрыты звёздочками

# Как скрыть отображение секретных параметров?

- В Spark есть специальный конфиг  
`spark.redaction.regex`
- Все параметры, названия которых соответствуют регулярке, будут скрыты звёздочками
- Или нет?

# Что я сделала

- Прокинула параметр

```
spark.yt.token=token
```

- Сказала спарку, что он секретный

```
spark.redaction.regex=(?i) secret|password|token
```



# Что получилось

- ✓ Параметр скрыт в SparkUI

# Но в логах драйвера...

- Запуск в Spark Standalone логирует строку запуска драйвера полностью

```
Launch Command: "java" "-cp"  
"-Dspark.yt.token=token" ...  
org.apache.spark.deploy.worker.DriverWrapper ...
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✗ Параметр виден в логах драйвера

# Что делать?

- В Spark 3.0 это починили

```
Launch Command: "java" "-cp"  
"-Dspark.yt.token=***** (redacted) " ...  
org.apache.spark.deploy.worker.DriverWrapper ...
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера

# Но в списке процессов...

- На хосте, где запущен spark-submit, видны все его аргументы

```
> ps aux | grep spark
```

```
spark-submit ... --conf spark.yt.token=token ...
```

# Может, обойти spark-submit?

- Всё равно придём к запуску класса SparkSubmit

```
> ps aux | grep spark
```

```
java -cp ... org.apache.spark.deploy.SparkSubmit  
  --conf spark.yt.token=token ...
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✗ Параметр виден в ps на хосте запуска spark-submit



# Что делать?

- Костыль в SparkSubmit.scala

```
sys.env.get("SPARK_YT_TOKEN").foreach {  
    token =>  
        sparkConf.set("spark.yt.token", token)  
}
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✓ Параметр скрыт в ps на хосте запуска spark-submit

# Но в списке процессов в PySpark...

- Пользователь PySpark в Jupyter заметил в своём списке процессов какую-то джаву с токеном в открытом виде

```
> ps aux | grep spark
```

```
python ...
```

```
java -cp ... org.apache.spark.deploy.SparkSubmit  
--conf spark.yt.token=token ...
```

# Прокидывание секретов

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✓ Параметр скрыт в ps на хосте запуска spark-submit
- ✗ Параметр виден в ps при запуске PySpark

# Что произошло?

- В коде для питона я прокидывала параметр при старте сессии

```
// example.py

spark = SparkSession.builder
    .config("spark.yt.token", token)
    .getOrCreate()
```

# Что произошло?

- В коде для питона я прокидывала параметр при старте сессии
- Но в PySpark вызов `getOrCreate` запускает jvm с `SparkSubmit`

```
// example.py
```

```
spark = SparkSession.builder  
    .config("spark.yt.token", token)  
    .getOrCreate()
```

# Что делать?

- У нас уже есть костыль в SparkSubmit, можем им воспользоваться

```
// example.py
```

```
os.environ["SPARK_YT_TOKEN"] = token  
spark = SparkSession.builder.getOrCreate()
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✓ Параметр скрыт в ps на хосте запуска spark-submit
- ✓ Параметр скрыт в ps при запуске PySpark



# Но в кластере на воркере...

- В списке процессов воркера, на котором запущен драйвер, видны все аргументы драйвера

```
> ps aux | grep spark
```

```
java -cp ... -Dspark.yt.token=token ...  
org.apache.spark.deploy.worker.DriverWrapper ...
```

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✓ Параметр скрыт в ps на хосте запуска spark-submit
- ✓ Параметр скрыт в ps при запуске PySpark
- ✗ Параметр виден в ps на воркере в кластере

# Что делать?

- Я добавила костыли
  - В запуск драйвера в SparkStandalone
  - В инициализацию SparkConf

# Что получилось

- ✓ Параметр скрыт в SparkUI
- ✓ Параметр скрыт в логах драйвера
- ✓ Параметр скрыт в ps на хосте запуска spark-submit
- ✓ Параметр скрыт в ps при запуске PySpark
- ✓ Параметр скрыт в ps на воркере в кластере

# Что делать?

- Не прокидывать секреты через конфиг спарка. Никогда. Даже если кажется, что это просто
- Использовать специальные сервисы для хранения секретов

# Итоги

- Простое решение было сделано за 3 дня
- Рефакторинг на более сложное занял ещё несколько недель
- Чтение батчами ускорило запросы в 2-3 раза

# В следующей серии...

- Pushdown predicate
- Оптимизация записи
- Использование информации о сортированности данных
- Гибкая интеграция с шедулерам

Спасибо! Вопросы?